



# A uniformly accurate numerical method for a class of dissipative systems

Philippe Chartier, Mohammed Lemou, Léopold Trémant

## ► To cite this version:

Philippe Chartier, Mohammed Lemou, Léopold Trémant. A uniformly accurate numerical method for a class of dissipative systems. *Mathematics of Computation*, 2022, 91 (334), pp.843-869. 10.1090/mcom/3688 . hal-02619512v2

**HAL Id: hal-02619512**

**<https://inria.hal.science/hal-02619512v2>**

Submitted on 1 Apr 2021

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# A uniformly accurate numerical method for a class of dissipative systems

Philippe Chartier\*, Mohammed Lemou†, Léopold Trémant‡

## Abstract

We consider a class of ordinary differential equations mixing slow and fast variations with varying stiffness (from non-stiff to strongly dissipative). Such models appear for instance in population dynamics or propagation phenomena. We develop a multi-scale approach by splitting the equations into a micro part and a macro part, from which the original stiffness has been removed. We then show that both parts can be simulated numerically with *uniform* order of accuracy using *standard explicit* numerical schemes. As a result, solving the problem in its micro-macro formulation can be done with a cost and an accuracy *independent of the stiffness*. This work is also a preliminary step towards the application of such methods to hyperbolic partial differential equations and we will indeed demonstrate that our approach can be successfully applied to two discretized hyperbolic systems (with and without non-linearities), though with some ad-hoc regularization.

*AMS subject classification (2020):* 65L04, 34E13, 65L05, 65L20

*Keywords:* dissipative problem, multi-scale, micro-macro decomposition, uniform accuracy

## 1 Introduction

We are interested in problems of the form, for  $x^\varepsilon(t) \in \mathbb{R}^{d_x}$  and  $z^\varepsilon(t) \in \mathbb{R}^{d_z}$ ,

$$\begin{cases} \dot{x}^\varepsilon = a(x^\varepsilon, z^\varepsilon), & x^\varepsilon(0) = x_0, \\ \dot{z}^\varepsilon = -\frac{1}{\varepsilon}Az^\varepsilon + b(x^\varepsilon, z^\varepsilon), & z^\varepsilon(0) = z_0, \end{cases} \quad (1.1)$$

with  $\varepsilon \in (0, 1]$  a small parameter,  $A$  a diagonal positive matrix with integer coefficients, and where  $a, b$  are respectively the  $x$ -component and the  $z$ -component of an analytic map  $f$  which smoothly depends on  $\varepsilon$ . We look for a solution  $x^\varepsilon(t), z^\varepsilon(t)$ , defined for  $t \in [0, 1]$ , irrespectively of the value of  $\varepsilon$ . The exact value of the right bound of the interval of definition of the solution, here 1, is somehow arbitrary, as it can be rescaled by changing the value of  $\frac{1}{\varepsilon}\Lambda$ . In the limit when  $\varepsilon$  goes to zero, the problem becomes stiff on the considered interval:

---

\*Inria Rennes, IRMAR and ENS Rennes, Campus de Beaulieu, F-35170 Bruz, France. Philippe.Chartier@inria.fr

†CNRS, IRMAR and ENS Rennes, Campus de Beaulieu, F-35170 Bruz, France. Mohammed.Lemou@univ-rennes1.fr

‡Inria Rennes and IRMAR, Campus de Beaulieu, 35049 Rennes, France. Leopold.Tremant@inria.fr

in other words, the problem resorts to long-time integration as 1 becomes large compared to  $\varepsilon$ . In the sequel we shall more often write the equations in compact form as

$$\dot{u}^\varepsilon = -\frac{1}{\varepsilon}\Lambda u^\varepsilon + f(u^\varepsilon), \quad u^\varepsilon(0) = u_0, \quad (1.2)$$

where  $u = \begin{pmatrix} x \\ z \end{pmatrix}$ ,  $\Lambda = \begin{pmatrix} 0 & 0 \\ 0 & A \end{pmatrix}$  and  $f(u) = \begin{pmatrix} a(x, z) \\ b(x, z) \end{pmatrix}$ . We set  $d = d_x + d_z$  the dimension of  $u$  such that  $u \in \mathbb{R}^d$ . Note that the dimension of  $x^\varepsilon$  may be zero without impacting our results. In contrast, it should be emphasized that we do not address the case where the map  $u \mapsto f(u)$  is a differential operator and  $u$  lies in a functional space: the theory required for that situation is outside the scope of our theorems. Nonetheless, two of our examples are discretized hyperbolic partial differential equations (PDEs) for which the method is successfully applied, even though an additional specific treatment is required.

Problems of the form 1.2 recurrently appear in population dynamics (see [GHM94; AP96; SAAP00; CCS18]), where  $A$  accounts for migration (in space and/or age) and  $a$  and  $b$  account for both the demographic and inter-population dynamics. In this context, the factor  $1/\varepsilon$  accounts for the fact that the migration dynamics is quantifiably faster than other dynamics involved.

When solving this kind of system numerically, problems arise due to the large range of values that  $\varepsilon$  can take. To be more specific, the error for standard methods of order  $q > 1$  behave like

$$E_\varepsilon(\Delta t) \leq \min \left( C_q \frac{\Delta t^q}{\varepsilon^r}, C_s \Delta t^s \right),$$

for some positive constants  $C_q$  and  $C_s$  independent of  $\varepsilon$  and integers  $s \leq q$  and  $r \geq 0$ . This forces very small values of  $\Delta t$  in order to achieve some accuracy and causes the computational cost of the simulation to increase greatly, often prohibitively so. Additionally, the order is reduced to  $s$  in the sense that<sup>1</sup>

$$\sup_{\varepsilon \in (0,1]} E_\varepsilon(\Delta t) \leq C \Delta t^s. \quad (1.3)$$

This behaviour is documented for instance in [HW96, Section IV.15] or in [HR07]. In order to ensure a given error bound, one must either accept this order reduction (if  $s > 0$ ), as is done for asymptotic-preserving (AP) schemes [Jin99] by taking a modified time-step  $\widetilde{\Delta t} = \Delta t^{q/s}$ , or use an  $\varepsilon$ -dependent time-step  $\Delta t = \mathcal{O}(\varepsilon^{r/q})$ .

A common approach to circumvent this difficulty is to invoke the *center manifold theorem* (see [Vas63; Car82; Sak90]), which dictates the long-time behaviour of the system and presents useful characteristics for numerical simulations: the dimension of the system is reduced and the dynamics on the manifold is non-stiff. However, this approach does not allow to capture the *transient phase* of the solution, i.e. the solution in short time before it reaches the stable manifold. Insofar as one wishes to describe the system out of equilibrium, this is clearly unsatisfactory. Furthermore, even if the solution is exponentially (w.r.t. time) close to the manifold, the center manifold approximation is accurate up to a certain error  $\mathcal{O}(\varepsilon^n)$ , rendering it useless if  $\varepsilon$  is of the order of 1.

---

<sup>1</sup>In particular, the scheme cannot be any usual explicit scheme since it would require a stability condition of the form  $\Delta t/\varepsilon < C$  with  $C$  independent of  $\varepsilon$ .

The strategy developed in this paper is based on a *micro-macro* decomposition of the problem in combination with the use of standard  $q^{th}$ -order *exponential Runge-Kutta* methods. It aims at deriving an overall scheme with an error  $E_\varepsilon(\Delta t)$  that can be bounded from above independently of  $\varepsilon$ , that is to say

$$E_\varepsilon(\Delta t) \leq C \Delta t^q$$

for some positive constant  $C$  independent of  $\varepsilon$ . In order to construct the appropriate transformation of the original system, we first provide a systematic way to compute asymptotic models at any order in  $\varepsilon$  approaching the solution over the *whole interval of time*. We then use the defect of this approximation to compute the solution with usual explicit numerical schemes and *uniform* accuracy (i.e. the cost and error of the scheme must be independent of  $\varepsilon$ ). This approach automatically overcomes the challenges posed by both extremes  $\varepsilon \ll 1$  and  $\varepsilon \sim 1$ .

The aforementioned micro-macro decomposition is obtained by writing the solution  $u^\varepsilon$  of (1.2) as the following composition of maps

$$u^\varepsilon(t) = \Omega_{t/\varepsilon}^\varepsilon \circ \Gamma_t^\varepsilon \circ (\Omega_0^\varepsilon)^{-1}(u_0) \quad (1.4)$$

where  $(\tau, u) \in \mathbb{R}_+ \times \mathbb{R}^d \mapsto \Omega_\tau^\varepsilon(u) \in \mathbb{R}^d$  is a change of variable  $\varepsilon$ -close to the map  $(\tau, u) \mapsto e^{-\tau\Lambda}u$  and where  $(t, u) \in [0, T] \times \mathbb{R}^d \mapsto \Gamma_t^\varepsilon(u)$  is the flow associated to a *non-stiff* autonomous vector field  $u \mapsto F^\varepsilon(u)$ , yet to be defined. The formal maps  $\Omega^\varepsilon$  and  $F^\varepsilon$  are approached at an arbitrary order  $n \in \mathbb{N}$  by  $\Omega^{[n]}$  and  $F^{[n]}$  respectively such that the equality

$$u^\varepsilon(t) = \Omega_{t/\varepsilon}^{[n]}(v^{[n]}(t)) + w^{[n]}(t) \quad (1.5)$$

holds true, where  $v^{[n]}(t) = \Gamma_t^{[n]} \circ (\Omega_0^{[n]})^{-1}(u_0)$  and  $w^{[n]}$  are respectively called the *macro* component and the *micro* component. A crucial feature of this decomposition is that  $w^{[n]}$  remains of size  $\mathcal{O}(\varepsilon^{n+1})$ .

Now, the main contribution of this work is to prove that, using explicit exponential Runge-Kutta (ERK) schemes of order  $n+1$  (which can be found for instance in [HO05]), it is possible to approximate  $u^\varepsilon$  with *uniform accuracy* and at *uniform computational cost* with respect to  $\varepsilon$ . In other words, we prove that formula (1.3) holds with  $s = q = n+1$  and  $r = 0$ . More precisely, if  $(t_i)_{0 \leq i \leq N}$  is a time-step grid of mesh-size  $\Delta t$ , and if  $(v_i)$  and  $(w_i)$  are computed numerically by applying the ERK method to the micro-macro decomposition, then there exists  $C$  independent of  $\varepsilon$  such that ( $|\cdot|$  stands for the usual Euclidian norm)

$$\max_{0 \leq i \leq N} \left\{ |x^\varepsilon(t_i) - x_i| + \frac{1}{\varepsilon} |z^\varepsilon(t_i) - z_i| \right\} \leq C \Delta t^{n+1} \quad \text{with} \quad \begin{pmatrix} x_i \\ z_i \end{pmatrix} = \Omega_{t_i/\varepsilon}^{[n]}(v_i) + w_i.$$

We emphasize here the expected occurrence of the scaling factor  $1/\varepsilon$  accounts for the fact that  $z$  becomes of size  $\mathcal{O}(\varepsilon)$  after a time  $\mathcal{O}(\varepsilon \log(1/\varepsilon))$ . IMEX methods such as CNLF and SBDF (see [ARW95; ACM99; HS21]), which mix implicit and explicit parts are not the focus of the article, but their use is briefly discussed in Remark 2.9.

The present work is related to the recent paper [CCS16], where asymptotic expansions of the solution of (1.1) are constructed for the special case where  $A$  is the identity matrix. The theory developed therein is however of no relevance for the construction of micro-macro decompositions as it relies heavily on trees and associated elementary differentials which can

hardly be computed in practice. Our approach actually shares more similarities with the one introduced for highly-oscillatory problems in [CLMV19] and later modified to become amenable for actual computations at any order [CLMZ20]. As a matter of fact, the technical arguments that sustain decomposition (1.4) are essentially adapted from [CCMM15] in a way that will be fully explained in Section 3.

The rest of the paper is organized as follows. In Section 2, we show our method to construct a micro-macro problem up to any order, and state our main result, i.e that solving this micro-macro problem with ERK schemes generates uniform accuracy on  $u^\varepsilon$ . In Section 3, we give proofs of all the results from Section 2. In Section 4, we present some techniques to adapt our method to discretized hyperbolic PDEs. Namely, we study a relaxed conservation law and the telegraph equation, which can be respectively found for instance in [JX95] and [LM08]. In Section 5, we verify our theoretical result of uniform accuracy by successfully obtaining uniform convergence when numerically solving micro-macro problems obtained from a toy ODE and from the two aforementioned discretized PDEs.

## 2 Uniform accuracy from a decomposition

We start by considering the solution  $u$  of

$$\partial_t u^\varepsilon = -\frac{1}{\varepsilon} \Lambda u^\varepsilon + f(u^\varepsilon), \quad u^\varepsilon(0) = u_0 \in \mathbb{R}^d, \quad (2.1)$$

and write it as the composition of a *non-stiff* flow  $(t, u) \mapsto \Gamma_t^\varepsilon(u)$  with a change of variable  $(\tau, u) \mapsto \Omega_\tau^\varepsilon(u)$  with  $\tau \in \mathbb{R}_+$ ,

$$u^\varepsilon(t) = \Omega_{t/\varepsilon}^\varepsilon \circ \Gamma_t^\varepsilon \circ (\Omega_0^\varepsilon)^{-1}(u_0). \quad (2.2)$$

In order for our approach to be rigorous, we start by introducing some definitions and assumptions in Subsection 2.1. We then present a way to approach these maps at any rank  $n \in \mathbb{N}$  by  $\Gamma^{[n]}$  and  $\Omega^{[n]}$  in Subsection 2.2. This approximation is such that the error in (2.2) is of size  $\mathcal{O}(\varepsilon^{n+1})$ . In Subsection 2.3, we use this approximation to construct a micro-macro problem which can be solved numerically using standard IMEX schemes. This leads to our main result: reconstructing the solution  $u^\varepsilon$  of (2.1) from the numerical solution of the micro-macro problem yields an error *independent of  $\varepsilon$*  on  $u^\varepsilon$ . All proofs are delayed until Section 3.

### 2.1 Definitions and assumptions

Before proceeding, we must first state the assumptions on the vector field  $u \mapsto f(u)$  and the operator  $\Lambda$ .

**Assumption 2.1.** *The matrix  $\Lambda$  is diagonal with nonnegative integer eigenvalues, and these values are nondecreasing when following the diagonal. In other words,  $\Lambda = \text{Diag}(\lambda_1, \dots, \lambda_d)$  with  $(\lambda_i)_{1 \leq i \leq d} \in \mathbb{N}^d$  and  $\lambda_1 \leq \dots \leq \lambda_d$ .*

Thanks to this assumption, we write  $u = \begin{pmatrix} x \\ z \end{pmatrix}$ , with  $(x, z)$  such that  $\Lambda u = \begin{pmatrix} 0 \\ Az \end{pmatrix}$  for some  $A$  positive definite. The dimension of  $z$  may be zero without making our results invalid.

**Assumption 2.2.** Let us set  $d_x$  and  $d_z$  the dimensions of  $x$  and  $z$  respectively. There exists a compact set  $X_1 \subset \mathbb{R}^{d_x}$  and a radius  $\check{\rho} > 0$  such that for every  $x$  in  $X_1$ , the map  $u \in \mathbb{R}^d \mapsto f(u) \in \mathbb{R}^d$  can be developed as a Taylor series around  $\begin{pmatrix} x \\ 0 \end{pmatrix}$ , and the series converges with a radius not smaller than  $\check{\rho}$ .

It is therefore possible to naturally extend  $f$  to compact subsets of  $\mathbb{C}^d$  defined by

$$\mathcal{U}_\rho := \left\{ u \in \mathbb{C}^d; \exists x \in X_1, \left| u - \begin{pmatrix} x \\ 0_{d_z} \end{pmatrix} \right| \leq \rho \right\},$$

for all  $0 \leq \rho < \check{\rho}$  as it is represented by a Taylor series in  $u \in \mathbb{C}^d$  on these sets. Here  $|\cdot|$  is the natural extension of the Euclidian norm on  $\mathbb{R}^d$  to  $\mathbb{C}^d$ .

It may seem particularly restrictive to assume that the  $z$ -component of the solution  $u^\varepsilon$  of (1.2) stays in a neighborhood of 0, however this is somewhat ensured by the *center manifold theorem*. This theorem states that there exists a map  $x \in \mathbb{R}^{d_x} \mapsto \varepsilon h^\varepsilon(x) \in \mathbb{R}^{d_z}$  smooth in  $\varepsilon$  and  $x$ , such that the manifold  $\mathcal{M}$  defined by

$$\mathcal{M} = \left\{ (x, z) \in \mathbb{R}^{d_x} \times \mathbb{R}^{d_z} : z = \varepsilon h^\varepsilon(x) \right\}$$

is a stable invariant for (1.1). It also states that all solutions  $(x^\varepsilon, z^\varepsilon)$  of (1.1) converge towards it exponentially quickly, i.e. there exists  $\mu > 0$  independent of  $\varepsilon$  such that

$$|z^\varepsilon(t) - \varepsilon h^\varepsilon(x^\varepsilon(t))| \leq C e^{-\mu t/\varepsilon}. \quad (2.3)$$

This means that the growth of  $z^\varepsilon$  is bounded by that of  $x^\varepsilon$ , and that after a time  $t \geq \varepsilon \log(1/\varepsilon)$ ,  $z^\varepsilon(t)$  is of size  $\mathcal{O}(\varepsilon)$ . Therefore it is credible to assume that  $z^\varepsilon$  stays somewhat close to 0. This is translated into a final assumption.

**Assumption 2.3.** There exist two radii  $0 < \rho_0 \leq \rho_1 < \check{\rho}$  and a closed subset  $X_0 \subset X_1 \subset \mathbb{R}^{d_x}$  such that the initial condition  $u_0 \in \mathbb{C}^d$  satisfies

$$\min_{x \in X_0} \left| u_0 - \begin{pmatrix} x \\ 0_{d_z} \end{pmatrix} \right| \leq \rho_0,$$

and for all  $\varepsilon \in (0, 1]$ , Problem (2.1) is well-posed on  $[0, 1]$  with its solution  $u^\varepsilon$  in  $\mathcal{U}_{\rho_1}$ .

Note that this is different to assuming that the initial data  $(x_0, z_0)$  is close to the center manifold. Indeed, the size of the initial condition is supposed independent of  $\varepsilon$ , therefore the distance from  $z(0)$  to the center manifold is always  $\mathcal{O}(1)$ .

For  $\rho \in [0, \check{\rho} - \rho_1)$ , we define the sets

$$\mathcal{K}_\rho := \mathcal{U}_{\rho_1 + \rho} = \left\{ u \in \mathbb{C}^d; \exists x \in X_1, \left| u - \begin{pmatrix} x \\ 0 \end{pmatrix} \right| \leq \rho_1 + \rho \right\} \quad (2.4)$$

which help quantify the distance to the solution  $u^\varepsilon$ . By Assumption 2.3, the solution of (1.2) is in  $\mathcal{K}_0$  at all time.

**Definition 2.4.** We introduce some technical constants:

- (i) A radius  $0 < R < \frac{1}{2}(\check{\rho} - \rho_1)$

(ii) An arbitrary rank  $p$  and a positive constant  $M$  such that for all  $0 \leq \alpha, \beta \leq p+2$  and all  $\sigma \in [0, 6\|\Lambda\|]$ ,

$$\frac{\sigma^\beta}{\beta!} \left\| (\rho_1 + 2R)^\alpha \partial_u^\alpha f \right\| \leq M$$

Given a radius  $0 \leq \rho \leq 2R$  and a map  $(\tau, u) \in \mathbb{R}_+ \times \mathcal{K}_\rho \mapsto \psi_\tau(u)$ , we define the norm,

$$\|\psi\|_\rho := \sup_{(\tau, u) \in \mathbb{R}_+ \times \mathcal{K}_\rho} |\psi_\tau(u)|. \quad (2.5)$$

If the map is furthermore  $p$ -times continuously differentiable w.r.t.  $\tau$ , then we define

$$\|\psi\|_{\rho, p} := \max_{0 \leq \nu \leq p} \|\partial_\tau^\nu \psi\|_\rho. \quad (2.6)$$

## 2.2 Constructing the micro-macro problem

We assume that the vector field in (2.2) follows an autonomous vector field  $F^\varepsilon$ , i.e.

$$\frac{d}{dt} \Gamma_t^\varepsilon(u) = F^\varepsilon(\Gamma_t^\varepsilon(u)). \quad (2.7)$$

Injecting this and (2.2) into (2.1) and writing  $v_0 = (\Omega_0^\varepsilon)^{-1}(u_0)$

$$(\partial_\tau + \Lambda) \Omega_{t/\varepsilon}^\varepsilon(\Gamma_t^\varepsilon(v_0)) = \varepsilon \left( f \circ \Omega_{t/\varepsilon}^\varepsilon(\Gamma_t^\varepsilon(v_0)) - \partial_u \Omega_{t/\varepsilon}^\varepsilon(\Gamma_t^\varepsilon(v_0)) \cdot F^\varepsilon(\Gamma_t^\varepsilon(v_0)) \right)$$

which by separation of scales  $t$  and  $t/\varepsilon$  generates the homological equation on  $\Omega^\varepsilon$ , for all  $(\tau, u) \in \mathbb{R}_+ \times K_\rho$ ,

$$(\partial_\tau + \Lambda) \Omega_\tau^\varepsilon(u) = \varepsilon (f \circ \Omega_\tau^\varepsilon(u) - \partial_u \Omega_\tau^\varepsilon(u) \cdot F^\varepsilon(u)). \quad (2.8)$$

It is furthermore possible to extract the vector field  $F^\varepsilon$  from this equation to get

$$F^\varepsilon = \langle \partial_u \Omega^\varepsilon \rangle^{-1} \langle f \circ \Omega^\varepsilon \rangle \quad (2.9)$$

where  $\langle \cdot \rangle$  is defined by the following formula

$$\langle \psi \rangle := \frac{1}{2\pi} \int_0^{2\pi} e^{i\theta\Lambda} \psi_{i\theta} d\theta, \quad (2.10)$$

with the canonical definition  $\psi_{i\theta} = \sum_{k \geq 0} e^{-ik\theta} \widehat{\psi}_k$ . To see this, we first observe that for an exponential series  $\tau \in \mathbb{R}_+ \mapsto \psi_\tau$  which converges absolutely for  $\tau = 0$ , i.e.  $\psi_\tau = \sum_{k \geq 0} e^{-k\tau} \widehat{\psi}_k$  with  $\sum_k \widehat{\psi}_k$  absolutely converging, we can extract the coefficient  $\widehat{\psi}_k$  as the Fourier coefficient of  $\psi_{i\theta}$  according to

$$\widehat{\psi}_k = \frac{1}{2\pi} \int_0^{2\pi} e^{ik\theta} \psi_{i\theta} d\theta. \quad (2.11)$$

Therefore, we write equation (2.8) as follows

$$\partial_\tau (e^{\tau\Lambda} \Omega_\tau^\varepsilon(u)) = \varepsilon (e^{\tau\Lambda} f \circ \Omega_\tau^\varepsilon(u) - e^{\tau\Lambda} \partial_u \Omega_\tau^\varepsilon(u) \cdot F^\varepsilon(u)), \quad (2.12)$$

and apply the Fourier operator (2.11) to get

$$\widehat{\partial_\tau(e^{\tau\Lambda}\Omega_\tau^\varepsilon)}(u)_k = \varepsilon \left( \widehat{(e^{\tau\Lambda}f \circ \Omega_\tau^\varepsilon(u))}_k - \widehat{(\partial_u \Omega_\tau^\varepsilon(u) \cdot F^\varepsilon(u))}_k \right).$$

Taking now  $k = 0$  and using definition (2.10) we get the expression (2.9). This framework of exponential series comes naturally thanks to Assumption 2.1.

The homological equation (2.8) has no unique solution in general, however we can approximate a solution as a *formal* solution as a power series in  $\varepsilon$ . This is generally the idea behind *normal forms*, where different methods have been developed (see [Mur06] for instance). Here we only consider a basic method to compute approximations  $\Omega^{[n]}$  and  $F^{[n]}$  of  $\Omega^\varepsilon$  and  $F^\varepsilon$  at any rank  $n \in \mathbb{N}$  by setting

$$(\partial_\tau + \Lambda)\Omega_\tau^{[n+1]} = \varepsilon(f \circ \Omega_\tau^{[n]} - \partial_u \Omega_\tau^{[n]} \cdot F^{[n]}). \quad (2.13)$$

with initial condition  $\Omega_\tau^{[0]} = e^{-\tau\Lambda}$ . Because we want  $\Omega^{[n+1]}$  to be an exponential series, it appears that necessarily,

$$F^{[n]} = \langle \partial_u \Omega^{[n]} \rangle^{-1} \langle f \circ \Omega^{[n]} \rangle. \quad (2.14)$$

However these equations alone are not enough to obtain  $\Omega^{[n]}$  at any order. Indeed, from (2.13), one gets

$$\Omega_\tau^{[n+1]} = e^{-\tau\Lambda}\Omega_0^{[n+1]} + \varepsilon \int_0^\tau e^{(\sigma-\tau)\Lambda} (f \circ \Omega_\sigma^{[n]} - \partial_u \Omega_\sigma^{[n]} \cdot F^{[n]}) d\sigma \quad (2.15)$$

meaning a choice of initial data  $\Omega_0^{[n+1]}$  is needed. One could think that choosing  $\Omega_0^{[n+1]} = \text{id}$  is the easiest choice, but computing (2.14) requires an inversion of  $\langle \partial_u \Omega^\varepsilon \rangle$ . Therefore we choose  $\Omega_0^{[n+1]}$  such that  $\langle \Omega^{[n+1]} \rangle = \text{id}$ , i.e. for all  $n \in \mathbb{N}$ ,

$$\Omega_0^{[n+1]} = \text{id} - \varepsilon \left\langle \int_0^\bullet e^{(\sigma-\bullet)\Lambda} (f \circ \Omega_\sigma^{[n]} - \partial_u \Omega_\sigma^{[n]} \cdot F^{[n]}) d\sigma \right\rangle \quad \text{thus} \quad F^{[n]} = \langle f \circ \Omega^{[n]} \rangle. \quad (2.16)$$

Now that we have a way to compute an approximate solution of (2.8), we introduce the error of approximation

$$\eta_\tau^{[n]} = \frac{1}{\varepsilon} (\partial_\tau + \Lambda) \Omega_\tau^{[n]} + \partial_u \Omega_\tau^{[n]} \cdot F^{[n]} - f \circ \Omega_\tau^{[n]}. \quad (2.17)$$

With these definitions, the maps  $(\tau, u) \mapsto \Omega_\tau^{[n]}(u)$ ,  $u \mapsto F^{[n]}(u)$  and  $(\tau, u) \mapsto \eta_\tau^{[n]}$  have the following properties.

**Theorem 2.5.** *For  $n$  in  $\mathbb{N}$ , let us denote  $r_n = R/(n+1)$  and  $\varepsilon_n := r_n/16M$  with  $R$  and  $M$  from Definition 2.4. For all  $\varepsilon > 0$  such that  $\varepsilon \leq \varepsilon_n$ , the maps  $(\tau, u) \mapsto \Omega_\tau^{[n]}(u)$ ,  $u \mapsto F^{[n]}(u)$  and  $(\tau, u) \mapsto \eta_\tau^{[n]}(u)$  given by (2.15) and (2.16) are well-defined on  $\mathbb{R}_+ \times \mathcal{K}_R$  and are analytic w.r.t.  $u$ . The change of variable  $\Omega^{[n]}$  and the residue  $\eta^{[n]}$  are both  $p+1$ -times continuously differentiable w.r.t.  $\tau$ . Moreover, with  $\|\cdot\|_R$  and  $\|\cdot\|_{R,p+1}$  given by (2.5) and (2.6), the following bounds are satisfied for all  $0 \leq \nu \leq p+1$ ,*

$$\begin{aligned} (i) \quad & \left\| \Omega^{[n]} - e^{-\tau\Lambda} \right\|_R \leq 4\varepsilon M, & (ii) \quad & \left\| \partial_\theta^\nu [\Omega^{[n]} - e^{-\tau\Lambda}] \right\|_R \leq 8(1 + \|\Lambda\|)^\nu \varepsilon M \nu! \\ (iii) \quad & \|F^{[n]}\|_R \leq 2M & (iv) \quad & \|\eta_\tau^{[n]}(u)\|_{R,p} \leq 2M(1 + \|\Lambda\|)^p \left( 2\mathcal{Q}_p \frac{\varepsilon}{\varepsilon_n} \right)^n \end{aligned}$$

where  $\|\cdot\|$  is the induced norm from  $\mathbb{R}^d$  to  $\mathbb{R}^d$ , and  $\mathcal{Q}_p$  is a  $p$ -dependent constant.

The proof will be treated in Subsection 3.1, and this results remains valid with the choice  $\Omega_0^{[n]} = \text{id}$ .



### 2.3 A result of uniform accuracy

Given a rank  $n \in \mathbb{N}$ , we now denote  $v^{[n]}(t) := \Gamma_t^{[n]} \circ (\Omega_0^{[n]})^{-1}(u_0)$  and inject the decomposition

$$u^\varepsilon(t) = \Omega_{t/\varepsilon}^{[n]}(v^{[n]}(t)) + w^{[n]}(t) \quad (2.18)$$

into Problem (2.1) in order to find an equation on  $w^{[n]}$ . The main interests of this decomposition can be roughly summarized as follows. First, the change of variable  $\Omega_{t/\varepsilon}^{[n]}$  is known explicitly and the macro solution  $v^{[n]}$  is smooth in  $\varepsilon$ , in the sense that time derivatives of  $v^{[n]}$  at any order are uniformly bounded with respect to  $\varepsilon \in (0, 1]$ . Second, the micro part  $w^{[n]}$  is less stiff than the original solution  $u^\varepsilon$  in the sense that its time derivatives, up to order  $n + 1$ , are uniformly bounded in  $\varepsilon$ . These important properties naturally allow the construction of numerical schemes on  $v^{[n]}$  and  $w^{[n]}$  that enjoy the *uniform accuracy*, i.e. in which the order of the numerical methods is independent of  $\varepsilon$  and is not degraded by the stiffness generated by the possibly small values of  $\varepsilon$ .

From decomposition (2.18) we obtain the following system

$$\begin{cases} \partial_t v^{[n]}(t) = F^{[n]}(v^{[n]}), \\ \partial_t w^{[n]}(t) = -\frac{1}{\varepsilon} \Lambda \left( \Omega_{t/\varepsilon}^{[n]}(v^{[n]}) + w^{[n]} \right) + f \left( \Omega_{t/\varepsilon}^{[n]}(v^{[n]}) + w^{[n]} \right) - \frac{d}{dt} \Omega_{t/\varepsilon}^{[n]}(v^{[n]}), \end{cases}$$

with initial conditions  $v^{[n]}(0) = (\Omega_0^{[n]})^{-1}(u_0)$  and  $w^{[n]}(0) = 0$ . By definition of  $v^{[n]}$  and using (2.17),

$$\begin{aligned} \frac{d}{dt} \Omega_{t/\varepsilon}^{[n]}(v^{[n]}(t)) &= \frac{1}{\varepsilon} \partial_\tau \Omega_{t/\varepsilon}^{[n]}(v^{[n]}) + \partial_u \Omega_{t/\varepsilon}^{[n]}(v^{[n]}) \cdot F^{[n]}(v^{[n]}) \\ &= -\frac{1}{\varepsilon} \Lambda \Omega_{t/\varepsilon}^{[n]}(v^{[n]}) + \eta_{t/\varepsilon}^{[n]}(v^{[n]}) + f(\Omega_{t/\varepsilon}^{[n]}(v^{[n]})). \end{aligned}$$

We get the micro-macro problem

$$\begin{cases} \partial_t v^{[n]}(t) = F^{[n]}(v^{[n]}), & (2.19a) \\ \partial_t w^{[n]}(t) = -\frac{1}{\varepsilon} \Lambda w^{[n]} + f \left( \Omega_{t/\varepsilon}^{[n]}(v^{[n]}) + w^{[n]} \right) - f \left( \Omega_{t/\varepsilon}^{[n]}(v^{[n]}) \right) - \eta_{t/\varepsilon}^{[n]}(v^{[n]}). & (2.19b) \end{cases}$$

with initial conditions  $v^{[n]}(0) = (\Omega_0^{[n]})^{-1}(u_0)$ ,  $w^{[n]}(0) = 0$ . The properties of this micro-macro problem can be summed up as followed.

**Theorem 2.6.** *For all  $n \in \mathbb{N}^*$ , let us define  $r_n = R/n$  and  $\varepsilon_n := r_n/16M$ , with  $R$  and  $M$  from Definition 2.4. For all  $\varepsilon \leq \varepsilon_n$ , Problem (2.19) is well-posed until some final time  $T_n$  independent of  $\varepsilon$ , and the following bounds are satisfied for all  $t \in [0, T_n]$  and  $0 \leq \nu \leq \min(n, p)$ ,*

$$\begin{aligned} (i) \quad & v^{[n]}(t) \in \mathcal{K}_R & (ii) \quad & |w^{[n]}(t)| \leq \frac{R}{4} \left( \frac{\varepsilon}{\varepsilon_n} \right)^{n+1} \\ (iii) \quad & |\partial_t^\nu E^{[n]}(t)| = \mathcal{O}(\varepsilon^{n-\nu}) & (iv) \quad & \|\partial_t^{\nu+1} E^{[n]}\|_{L^1} = \mathcal{O}(\varepsilon^{n-\nu}) \end{aligned}$$

where  $E^{[n]} = \partial_t w^{[n]} + \frac{1}{\varepsilon} \Lambda w^{[n]}$ .

**Remark 2.7.** *The attentive reader may notice that, while we made the computation of  $F^{[n]}$  easy with (2.16), the initial condition of the macro part,  $v^{[n]}(0) = (\Omega_0^{[n]})^{-1}(u_0)$ , is not explicit. However, this system must be solved only once, while  $F^{[n]}$  is used at every time-step. Furthermore, it is possible to compute an approximation of  $v^{[n]}(0)$  explicitly up to  $\mathcal{O}(\varepsilon^{n+1})$  using<sup>2</sup>*

$$v^{[n+1]}(0) = u_0 - \left(\Omega_0^{[n+1]} - \text{id}\right)(v^{[n]}(0)) + \mathcal{O}(\varepsilon^{n+2}) \quad (2.20)$$

with initialization  $v^{[0]}(0) = u_0$ . Because  $\Omega_0^{[n+1]}$  is near-identity (up to  $\mathcal{O}(\varepsilon)$ ), an error of size  $\varepsilon^{n+1}$  on  $v^{[n]}(0)$  will only translate in an error of size  $\varepsilon^{n+2}$  on  $v^{[n+1]}(0)$ .

We can now define approached initial conditions for the micro-macro problem iterating (2.20) at each rank  $n$  and truncating the  $\mathcal{O}(\varepsilon^{n+2})$  term. The initial condition of the micro part becomes

$$w^{[n]}(0) = u_0 - \Omega_0^{[n]}(v_n) \quad (2.21)$$

which ensures  $w^{[n]}(0) = \mathcal{O}(\varepsilon^{n+1})$ , meaning our results are not jeopardised.

Using a standard explicit scheme to solve Problem (2.19) cannot work due to the term  $\frac{1}{\varepsilon}\Lambda w^{[n]}$ . This is why we focus on exponential schemes, which render this term non-problematic in terms of stability (see [MZ09]). Of course, the only use of these exponential schemes does not solve the problem of non-uniform order of accuracy however, as these schemes all reduce to order 1 when taking the supremum of the error for  $\varepsilon \in (0, \varepsilon^*]$ . This is where our micr-macro formulation plays a crucial role since it allows standard numerical schemes (like exponential Runge-Kutta schemes for instance) to *keep their order uniformly* in  $\varepsilon \in (0, 1]$ . It should be noted that exponential schemes are well-established and the formulas to implement them can be found for example in [HO05] up to the fourth-order.

The first-order Euler method applied to (1.2) would yield

$$u_{i+1} = e^{-\frac{\Delta t}{\varepsilon}\Lambda} u_i + \Delta t \varphi\left(-\frac{\Delta t}{\varepsilon}\Lambda\right) f(u_i)$$

with  $\varphi(-h\Lambda) = \frac{1}{h} \int_0^h e^{-s\Lambda} ds$ . Because  $\Lambda$  is diagonal, this type of integral is easy to compute. There is no computational drawback to exponential schemes in this case. Furthermore, for these schemes the error bound involves the "modified" norm

$$|u|_\varepsilon = \left| u + \frac{1}{\varepsilon}\Lambda u \right|. \quad (2.22)$$

This norm is interesting because after a short time  $t \geq \varepsilon \log(1/\varepsilon)$ , the  $z$ -component of the solution  $u^\varepsilon$  of (1.2) is of size  $\varepsilon$ , as evidenced by the center manifold theorem in (2.3). Using the norm  $|\cdot|_\varepsilon$  somewhat rescales  $z^\varepsilon$  (but not  $x^\varepsilon$ ) by  $\varepsilon^{-1}$  such that studying the error in this norm can be seen as a sort of "relative" error.

The following result asserts that, indeed, our micro-macro reformulation of the problem allows any numerical scheme of order  $p$ , namely exponential schemes, to enjoy the uniform accuracy property, with the same order  $p$ . A detailed presentation of exponential Runge-Kutta schemes can be found for instance in [HO05; HO04].

---

<sup>2</sup>The above formula is a consequence of the behaviour of the error,  $\Omega^{[n+1]} = \Omega^{[n]} + \mathcal{O}(\varepsilon^{n+1})$  (see [CLMV19]), therefore  $v^{[n+1]}(0) = v^{[n]}(0) + \mathcal{O}(\varepsilon^{n+1})$ . Injecting this last approximation in  $v^{[n+1]}(0) = u_0 - (\Omega^{[n+1]} - \text{id})(v^{[n+1]}(0))$  generates the formula.

**Theorem 2.8.** *Under the assumptions of Theorem 2.6 and denoting  $T_n \leq T$  a final time such that Problem (2.19) is well-posed on  $[0, T_n]$ . Given  $(t_i)_{i \in \llbracket 0, N \rrbracket}$  a discretisation of  $[0, T_n]$  of time-step  $\Delta t := \max_i |t_{i+1} - t_i|$ . computing an approximate solution  $(v_i, w_i)$  of (2.19) using an exponential Runge-Kutta scheme of order  $q := \min(n, p) + 1$  yields a uniform error of order  $q$ , i.e.*

$$\max_{0 \leq i \leq N} |u^\varepsilon(t_i) - \Omega_{t_i/\varepsilon}^{[n]}(v_i) - w_i|_\varepsilon \leq C \Delta t^q \quad (2.23)$$

where  $C$  is independent of  $\varepsilon$ .

The left-hand side of this inequality involves  $|\cdot|_\varepsilon$  and shall be called the modified error. It dominates the absolute error which uses  $|\cdot|$ .

**Remark 2.9.** *Only exponential schemes are considered here rather than for instance IMEX-BDF schemes which are sometimes preferred (as in [HS21]). The reason for this is twofold.*

*First, as was mentioned already, iterations are easy to compute because of the diagonal nature of  $\Lambda$ . Second, the error bounds are generally better for these schemes. Indeed, an IMEX-BDF scheme of order  $q$  involves the  $L^1$  norm of  $\partial_t^{q+1} w^{[n]}$ , which is worse than the  $L^1$  norm of  $\partial_t^q E^{[n]}$ . The former is of size  $\mathcal{O}(\varepsilon^{n-q})$  while the latter is of size  $\mathcal{O}(\varepsilon^{n+1-q})$ . We made the choice to prioritize methods of order  $n+1$  rather than  $n$ .*

### 3 Proofs of theorems from Section 2

#### 3.1 Proof of Theorem 2.5: properties of the decomposition

For some rank  $n \in \mathbb{N}$ , consider the change of variable  $(\tau, u) \mapsto \Omega_\tau^{[n]}(u)$  given by (2.15) and (2.16). From a straightforward induction using Assumptions 2.1 and 2.2, it appears that this change of variable can be written as a *formal* exponential series,

$$\Omega_\tau^{[n]}(u) = \sum_{k \in \mathbb{N}} e^{-k\tau} \widehat{\Omega^{[n]}_k}(u).$$

This can be associated to a power series  $\Xi^{[n]}(\xi; u) = \sum_{k \in \mathbb{N}} \xi^k \widehat{\Omega^{[n]}_k}(u)$ ,  $\xi \in \mathbb{C}$ ,  $|\xi| \leq 1$ , which is entirely determined by its behaviour on the border, i.e. by the periodic map

$$\Phi_\theta^{[n]}(u) = \Xi^{[n]}(e^{i\theta}; u) = \Omega_{-i\theta}^{[n]}(u) = \sum_{k \in \mathbb{N}} e^{ik\theta} \widehat{\Omega^{[n]}_k}(u). \quad (3.1)$$

Differentiating  $\Phi^{[n+1]}$  w.r.t.  $\theta$  and identifying the coefficients in (2.13), we obtain a (still formal) homological equation on  $\Phi^{[n]}$ :

$$(\partial_\theta - i\Lambda) \Phi_\theta^{[n+1]} = -i\varepsilon (f \circ \Phi_\theta^{[n]} - \partial_u \Phi_\theta^{[n]} \cdot F^{[n]}). \quad (3.2)$$

The periodic defect  $\delta_\theta^{[n]} = -i\eta_{-i\theta}^{[n]}$  satisfies

$$\delta_\theta^{[n]} = \frac{1}{\varepsilon} (\partial_\theta - i\Lambda) \Phi_\theta^{[n]} + i f \circ \Phi_\theta^{[n]} - i \partial_u \Phi_\theta^{[n]} \cdot F^{[n]} \quad (3.3)$$

Note that these relations both use the identity

$$\sum_{k \in \mathbb{N}} \xi^k \widehat{f \circ \Omega^{[n]}_k} = f \left( \sum_{k \in \mathbb{N}} \xi^k \widehat{\Omega^{[n]}_k} \right) \quad (3.4)$$

which seems fairly evident, but requires the right-hand side of the equation to be well-defined for all  $|\xi| \leq 1$ .

Setting the filtered map  $\tilde{\Phi}_\theta^{[n]} = e^{-i\theta\Lambda}\Phi_\theta^{[n]}$ , it satisfies

$$\partial_\theta \tilde{\Phi}_\theta^{[n+1]} = \varepsilon \left( g_\theta \circ \tilde{\Phi}_\theta^{[n]} - \partial_u \tilde{\Phi}_\theta^{[n]} \cdot G^{[n]} \right) \quad (3.5)$$

with  $g_\theta(u) = e^{-i\theta\Lambda}f(e^{i\theta\Lambda}u)$  and  $G^{[n]} = iF^{[n]}$ .

**Property 3.1.** *Assumptions 2.2 and 2.3 ensure the following properties, with  $R, M$  and  $p$  given in Definition 2.4:*

- (i) *For all  $\varepsilon \in (0, 1]$ , the Cauchy problem  $\partial_t y^\varepsilon = g_{t/\varepsilon}(y^\varepsilon)$ ,  $y^\varepsilon(0) = u_0$  is well-posed in  $\mathcal{K}_0$  up to some final time independent of  $\varepsilon$ .*
- (ii) *For all  $\theta \in \mathbb{T}$ , the function  $u \mapsto g_\theta(u)$  is analytic from  $\mathcal{K}_{2R}$  to  $\mathbb{C}^d$ .*
- (iii) *For all  $\sigma \in [0, 3]$ ,*

$$\forall 0 \leq \nu \leq p+2, \quad \frac{\sigma^\nu}{\nu!} \|\partial_\theta^\nu g\|_{\mathbb{T}, 2R} \leq M, \quad (3.6)$$

Initial condition (2.16) means that the periodic change of variable would be defined by

$$\tilde{\Phi}_\theta^{[n+1]} = \text{id} + \varepsilon(T_\theta^{[n]} - \Pi(T^{[n]})) \quad \text{and} \quad \Phi_\theta^{[n+1]} = e^{i\theta\Lambda}\tilde{\Phi}_\theta^{[n+1]} \quad (3.7)$$

with  $\Pi$  the average<sup>3</sup> and  $T_\theta^{[n]} = \int_0^\theta (g_\sigma \circ \tilde{\Phi}_\sigma^{[n]} - \partial_u \tilde{\Phi}_\sigma^{[n]} \cdot G^{[n]}) d\sigma$ . Because  $\tilde{\Phi}^{[n]}$  is periodic at all rank  $n$ , taking the average in (3.5) gives the vector field

$$G^{[n]} = \Pi(g \circ \tilde{\Phi}^{[n]}). \quad (3.8)$$

This is known as *standard averaging*. We introduce norms on periodic maps akin to (2.5) and (2.6), namely for  $0 \leq \rho \leq 2R$ , given a periodic map  $(\theta, u) \in \mathbb{T} \times \mathcal{K}_\rho \mapsto \varphi_\theta(u)$ ,

$$\|\varphi\|_{\mathbb{T}, \rho} := \sup_{(\theta, u) \in \mathbb{T} \times \mathcal{K}_\rho} |\varphi_\theta(u)| \quad \text{and} \quad \|\varphi\|_{\mathbb{T}, \rho, \nu} := \max_{0 \leq \alpha \leq \nu} \|\varphi_\theta(u)\|_{\mathbb{T}, \rho} \quad (3.9)$$

where the second norm assumes that  $\varphi$  is  $\nu$ -times continuously differentiable w.r.t.  $\theta$ . Then the following bounds are satisfied.

**Theorem 3.2** (from [CLMV19] and [CCMM15]). *For  $n \in \mathbb{N}$ , let us denote  $r_n = R/(n+1)$  and  $\varepsilon_n := r_n/16M$ . For all  $\varepsilon > 0$  such that  $\varepsilon \leq \varepsilon_n$ , the maps  $\Phi^{[n]}$  and  $G^{[n]}$  are well-defined by (3.7) and (3.8). The change of variable  $\Phi^{[n]}$  and the defect  $\delta^{[n]}$  are both  $(p+2)$ -times continuously differentiable w.r.t.  $\theta$ , and  $\Phi_0^{[n]}$  is invertible with analytic inverse on  $\mathcal{K}_{R/4}$ . Moreover, the following bounds are satisfied for  $1 \leq \nu \leq p+1$ ,*

$$\begin{aligned} (i) \quad & \|\tilde{\Phi}^{[n]} - \text{id}\|_{\mathbb{T}, R} \leq 4\varepsilon M \leq \frac{r_n}{4}, & (ii) \quad & \|\partial_\theta^\nu \tilde{\Phi}^{[n]}\|_{\mathbb{T}, R} \leq 8\varepsilon M \nu! \\ (iii) \quad & \|G^{[n]}\|_{\mathbb{T}, R} \leq 2M & (iv) \quad & \|\tilde{\delta}^{[n]}\|_{\mathbb{T}, R, p+1} \leq 2M \left( 2\mathcal{Q}_p \frac{\varepsilon}{\varepsilon_n} \right)^n \end{aligned}$$

where  $\tilde{\Phi}_\theta^{[n]} = e^{-i\theta\Lambda}\Phi_\theta^{[n]}$  and  $\tilde{\delta}^{[n]} = e^{-i\theta\Lambda}\delta_\theta^{[n]}$  correspond to the filtered equation (3.5), and  $\mathcal{Q}_p$  is a  $p$ -dependent constant.

---

<sup>3</sup>Explicitely,  $\Pi(\varphi) = \frac{1}{2\pi} \int_0^{2\pi} \varphi_\sigma d\sigma$

In order to prove Theorem 2.5, we show that the previous calculations of this section are rigorous rather than formal. Let us work by induction and assume that the negative modes of  $\Phi^{[n]}$  vanish (this is true for  $\Phi_\theta^{[0]} = e^{i\theta\Lambda}$  since  $\Lambda$  is positive semidefinite). Because  $(\theta, u) \mapsto \Phi_\theta^{[n]}(u)$  is continuously differentiable w.r.t.  $\theta$ , its Fourier series converges absolutely, thus  $(\xi, u) \mapsto \Xi^{[n]}(\xi; u)$  is well-defined for all  $|\xi| \leq 1$  and  $u \in \mathcal{K}_R$ . By maximum modulus principle,

$$\|\Omega^{[n]} - e^{-\tau\Lambda}\|_R \leq \sup_{|\xi| \leq 1, u \in \mathcal{K}_R} |\Xi^{[n]}(\xi; u) - \xi^\Lambda| \leq \|\Phi^{[n]} - e^{i\theta\Lambda}\|_{\mathbb{T}, R} \leq \|\tilde{\Phi}^{[n]} - \text{id}\|_{\mathbb{T}, R}$$

The reasoning also stands for all derivatives  $1 \leq \nu \leq p+1$ ,

$$\left\| \partial_\tau^\nu [\Omega^{[n]} - e^{-\tau\Lambda}] \right\| \leq \sup_{\xi, u} \left| (\xi \partial_\xi)^\nu [\Xi^{[n]}(\xi; u) - \xi^\Lambda] \right| \leq \left\| \partial_\theta^\nu [\Phi^{[n]} - e^{i\theta\Lambda}] \right\|_{\mathbb{T}, R}$$

and  $\left\| \partial_\theta^\nu [\Phi^{[n]} - e^{i\theta\Lambda}] \right\|_{\mathbb{T}, R} \leq (1 + \|\Lambda\|)^\nu \|\partial_\theta^\nu \tilde{\Phi}^{[n]}\|_{\mathbb{T}, R, \nu}$ . Furthermore, for  $u \in \mathcal{K}_R$ , let  $x \in X_1$  s.t.  $\left| u - \begin{pmatrix} x \\ 0 \end{pmatrix} \right| \leq \rho_1 + R$ . Then for all  $|\xi| \leq 1$ ,

$$\left| \Xi^{[n]}(\xi; u) - \begin{pmatrix} x \\ 0 \end{pmatrix} \right| \leq \left| \Phi_\theta^{[n]}(u) - \begin{pmatrix} x \\ 0 \end{pmatrix} \right| \leq \left| \tilde{\Phi}_\theta^{[n]}(u) - \begin{pmatrix} x \\ 0 \end{pmatrix} \right|,$$

since a multiplication by  $e^{-i\theta\Lambda}$  has no influence on the norm, nor on  $\begin{pmatrix} x \\ 0 \end{pmatrix}$ . A triangle inequality yields

$$\left| \Xi^{[n]}(\xi; u) - \begin{pmatrix} x \\ 0 \end{pmatrix} \right| \leq |\tilde{\Phi}^{[n]} - u| + \left| u - \begin{pmatrix} x \\ 0 \end{pmatrix} \right| < \rho_1 + 2R,$$

therefore  $f(\Xi^{[n]}(\xi; u))$  is well-defined for all  $|\xi| \leq 1$  and  $u \in \mathcal{K}_R$ , by expanding it into an absolutely converging series around  $\begin{pmatrix} x \\ 0 \end{pmatrix}$ , thereby justifying relations (3.4) and (3.3). The maximum modulus principle can finally be applied to the couple  $(\eta^{[n]}, \delta^{[n]})$  in order to obtain the last estimate of Theorem 2.5.  $\square$

### 3.2 Proof of Theorem 2.6: well-posedness of the micro-macro problem

This proof is in several parts: first we show that problem (2.19a) is well-posed, and use this result to show that the bound on  $w^{[n]}$  is satisfied, thereby also proving that (2.19b) is well-posed. Finally we focus on the bounds on  $E^{[n]}$ .

Let us set  $\varphi(v) = u_0 + v - \Omega_0^{[n]}(u_0 + v)$ . Using Theorem 2.5, if  $|v| \leq R/4$  then  $|\varphi(v)| \leq R/4$ . By Brouwer fixed-point theorem, there exists  $v^*$  such that  $\varphi(v^*) = v^*$ , i.e.  $u^* \in \mathcal{K}_{R/4}$  such that  $\Omega_0^{[n]}(u^*) = u_0$ . Therefore  $v^{[n]}(0) := u^* \in \mathcal{K}_{R/4}$ .

Given  $t > 0$  and assuming  $v^{[n]}(s) \in \mathcal{K}_R$  for all  $s \in [0, t]$ , one can bound  $v^{[n]}(t)$  using Theorem 2.5:

$$\left| v^{[n]}(t) - v^{[n]}(0) \right| = \left| \int_0^t F^{[n]}(v^{[n]}(s)) \, ds \right| \leq 2Mt.$$

Setting  $T_v := \frac{3R}{8M}$  ensures  $|v^{[n]}(t) - v^{[n]}(0)| \leq 3R/4$ , meaning that for all  $t \in [0, T_v]$ ,  $v^{[n]}(t)$  exists and is in  $\mathcal{K}_R$ . Again from Theorem 2.5, we deduce  $\Omega_\tau^{[n]}(v^{[n]}(t)) \in \mathcal{K}_{5R/4}$ .

Focusing now on  $w^{[n]}$  and assuming for all  $s \in [0, t]$ ,  $|w^{[n]}(s)| \leq R/4$ , the linear term  $L^{[n]}(\tau, s, w^{[n]}(s))$  is bounded using a Cauchy estimate:

$$|L^{[n]}(\tau, s, w^{[n]}(s))| \leq \|\partial_u f\|_{3R/2} \leq \frac{\|f\|_{2R}}{2R - \frac{3}{2}R} \leq \frac{2M}{R}$$

using a Cauchy estimate. The integral form then gives the bounds

$$\begin{aligned} |w^{[n]}(t)| &\leq \left| \int_0^t e^{\frac{s-t}{\varepsilon}\Lambda} L^{[n]}(s/\varepsilon, s, w^{[n]}(s)) w^{[n]}(s) ds + \int_0^t e^{\frac{s-t}{\varepsilon}\Lambda} S^{[n]}(s/\varepsilon, s) ds \right| \\ &\leq \int_0^t \frac{2M}{R} |w^{[n]}(s)| ds + \left| \int_0^t e^{\frac{s-t}{\varepsilon}\Lambda} S^{[n]}(s/\varepsilon, s) ds \right| \end{aligned} \quad (3.10)$$

Using the notation of the previous subsection,  $\tilde{\delta}_\theta^{[n]} = -ie^{-i\theta\Lambda}\eta_{-i\theta}^{[n]}$ , from which

$$\eta_\tau^{[n]}(u) = \sum_{k \in \mathbb{Z}} e^{-(k+\Lambda)\tau} c_k^{[n]}(u) \quad \text{with} \quad c_k^{[n]}(u) = \frac{1}{2\pi} \int_0^{2\pi} e^{-ik\theta} \tilde{\delta}_\theta^{[n]}(u) d\theta.$$

Since  $\langle \eta^{[n]} \rangle = 0$ , i.e.  $c_0^{[n]} = 0$ , it is possible to bound the source term in  $w^{[n]}$  by

$$\begin{aligned} \left| \int_0^t e^{\frac{s-t}{\varepsilon}\Lambda} S^{[n]}(s/\varepsilon, s) ds \right| &\leq \left\| e^{-\frac{t}{\varepsilon}\Lambda} \right\| \int_0^t \sum_{k \in \mathbb{Z}^*} (e^{-k\frac{s}{\varepsilon}} \|c_k^{[n]}\|_{\mathbb{T}, R}) ds \\ &\leq \sum_{k \in \mathbb{Z}^*} \frac{\varepsilon}{k} \|c_k^{[n]}\|_{\mathbb{T}, R} \leq \varepsilon \left( \sum_{k \in \mathbb{Z}^*} \frac{1}{k^2} \right) \|\partial_\theta \tilde{\delta}^{[n]}\|_{\mathbb{T}, R} \end{aligned}$$

where  $\|\cdot\|_{\mathbb{T}, R}$  is given by (3.9). Using Theorem 3.2, there exists a constant  $M_n > 0$  such that for all  $t \in [0, T_v]$ ,

$$\left| \int_0^t e^{\frac{s-t}{\varepsilon}\Lambda} S^{[n]}(s/\varepsilon, s) ds \right| \leq M_n \left( \frac{\varepsilon}{\varepsilon_n} \right)^{n+1}. \quad (3.11)$$

Using Gronwall's lemma in (3.10) with this inequality yields

$$|w^{[n]}(t)| \leq M_n e^{\frac{2M}{R}t} \left( \frac{\varepsilon}{\varepsilon_n} \right)^{n+1} \leq M_n e^{\frac{2M}{R}t}.$$

We now set  $T_w > 0$  such that  $M_n e^{\frac{2M}{R}T_w} \leq R/4$  ( $T_w$  may therefore depend on  $n$ , but does not depend on  $\varepsilon$ ) and

$$T_n = \min(T_v, T_w).$$

This ensures the well-posedness of (2.19) on  $[0, T_n]$  as well as the size of  $w^{[n]}$ .

Finally, the results on  $E^{[n]}$  are a direct consequence of the bounds on the linear term

$$\sup_{\alpha+\beta+\gamma \leq p+1} \|\partial_\tau^\alpha \partial_t^\beta \partial_u^\gamma L^{[n]}\| < +\infty$$

and on the source term

$$\sup_{0 \leq \alpha+\beta \leq p} \|\partial_\tau^\alpha \partial_t^\beta S^{[n]}\|_{L^\infty} = \mathcal{O}(\varepsilon^n), \quad \sup_{\substack{\beta \geq 1 \\ 1 \leq \alpha+\beta \leq p+1}} \|\partial_\tau^\alpha \partial_t^\beta S^{[n]}\|_{L^1} = \mathcal{O}(\varepsilon^{n+1}).$$

This stems directly from Cauchy estimates and Theorem 2.5. □

### 3.3 Proof of Theorem 2.8: uniform accuracy

The idea in this proof is to bound the errors on the macro part and micro part separately, using

$$\left| u^\varepsilon(t_i) - \Omega_{t_i/\varepsilon}^{[n]}(v_i) - w_i \right|_\varepsilon \leq \left| \Omega_{t_i/\varepsilon}^{[n]}(v^{[n]}(t_i)) - \Omega_{t_i/\varepsilon}^{[n]}(v_i) \right|_\varepsilon + \left| w^{[n]}(t_i) - w_i \right|_\varepsilon.$$

As the macro part  $v^{[n]}$  involves no linear term, the scheme acts like any RK scheme on this part. Since  $v^{[n]}$  and  $F^{[n]}$  are non-stiff, the scheme is necessarily *uniformly* of order  $q$ , i.e.

$$\left| v^{[n]}(t_i) - v_i \right| \leq \Delta t^q \cdot t_i \cdot \|\partial_t^{q+1} v^{[n]}\|_{L^\infty}$$

using usual error bounds on RK schemes. The reader may notice that the absolute error involving  $|\cdot|$  was used, not the modified error involving  $|\cdot|_\varepsilon$ . The results in [HO04] state that an exponential RK scheme of order  $q$  generates an error given by

$$\left| w^{[n]}(t_i) - w_i \right|_\varepsilon \leq C \Delta t^q \left( \|\partial_t^{q-1} E^{[n]}\|_\infty + \|\partial_t^q E^{[n]}\|_{L^1} \right). \quad (3.12)$$

The bounds on  $E^{[n]} = \partial_t w^{[n]} + \frac{1}{\varepsilon} \Lambda w^{[n]}$  and its derivatives w.r.t.  $\varepsilon$  can be found in Theorem 2.6, rendering the computation of bounds on the error of the micro part straightforward. From Theorem 2.5.(i),  $\Omega_\tau^{[n]}(u) = e^{-\tau\Lambda} u + \mathcal{O}(\varepsilon)$ , therefore the error on  $\Omega_{t/\varepsilon}^{[n]}(v^{[n]})$  is of the form

$$\Omega_{t_i/\varepsilon}^{[n]}(v^{[n]}(t_i)) - \Omega^{[n]}(v_i) = e^{-t_i\Lambda/\varepsilon} (v^{[n]}(t_i) - v_i) + \varepsilon r_i$$

where  $v^{[n]}(t_i) - v_i$  and  $r_i$  are of size  $t_i \cdot \Delta t^q$ . The error can therefore be bounded, denoting  $\|\cdot\|$  the induced norm from  $\mathbb{R}^d$  to  $\mathbb{R}^d$ ,

$$\left| \Omega_{t_i/\varepsilon}^{[n]}(v^{[n]}(t_i)) - \Omega^{[n]}(v_i) \right|_\varepsilon \leq \left( 1 + \left\| \frac{t_i}{\varepsilon} \Lambda e^{-\frac{t_i}{\varepsilon} \Lambda} \right\| \right) |v^{[n]}(t_i) - v_i| + (\varepsilon + \|\Lambda\|) |r_i|.$$

From this we get the desired result on  $u^\varepsilon$ . □

## 4 Application to some ODEs derived from discretized PDEs

In this section, we construct micro-macro problems for two *discretized* hyperbolic relaxation systems of the form

$$\begin{cases} \partial_t u + \partial_x \tilde{u} = 0 \\ \partial_t \tilde{u} + \partial_x u = \frac{1}{\varepsilon} (g(u) - \tilde{u}) \end{cases}$$

where  $g$  acts either as a differential operator on  $u$  (telegraph equation, Subsection 4.1), or as a scalar value function (relaxed conservation law, Subsection 4.2). These two problems may seem similar in theory, and the latter actually serves as a stepping stone to treat the former in [JPT98; JPT00], but we will treat them quite differently in practice. Some recent AP schemes with promising convergence have been developed for this type of problems in [BPR17; ADP20].

Let us insist that we only consider these problems *after discretization* (using either Fourier modes or an upwind scheme), yet even in a discrete framework, it will be apparent that a direct application of the method is impossible, often because of the apparition of a backwards heat equation. The goal of this section is precisely to present some possible workarounds to overcome the problems that appear. Should the reader wish to see a more detailed and direct application of our method, they can find one in Subsection 5.1.

## 4.1 The telegraph equation

A commonly studied equation in kinetic theory is the one-dimensional Goldstein-Taylor model, also known as the telegraph equation (see [JPT98; LM08], for instance). It can be written, for  $(t, x) \in [0, T] \times \mathbb{R}/2\pi\mathbb{Z}$

$$\begin{cases} \partial_t \rho + \partial_x j = 0, \\ \partial_t j + \frac{1}{\varepsilon} \partial_x \rho = -\frac{1}{\varepsilon} j, \end{cases} \quad (4.1)$$

where  $\rho$  and  $j$  represent the mass density and the flux respectively. Using Fourier transforms in  $x$ , it is possible to represent a function  $v(t, x)$  by

$$v(t, x) = \sum_{k \in \mathbb{Z}} v_k(t) e^{ikx}.$$

Considering a given frequency  $k \in \mathbb{Z}$  the problem can be reduced to

$$\begin{cases} \partial_t \rho_k = -ik j_k, \\ \partial_t j_k = -\frac{1}{\varepsilon} (j_k + ik \rho_k). \end{cases}$$

Treating this problem using our method directly leads to dead-ends, therefore we will guide the reader through our reasoning navigating some of these dead-ends. This will lead to micro-macro decompositions of orders 0 and 1. These struggles can be seen as limitations of our approach, however we show that with only slight tweaks, it is possible to obtain an error of uniform order 2 using a standard exponential RK scheme. This result is summed up at the end of this subsection as Proposition 4.1.

In order to make a component  $-\frac{1}{\varepsilon} z_k$  appear, it would be tempting to set  $z_k = j_k + ik \rho_k$ . This quantity would verify the following differential equation

$$\partial_t z_k = -\frac{1}{\varepsilon} z_k + k^2 z_k - ik^3 \rho_k.$$

Integrating this differential equation gives

$$z_k(t) = \exp\left(-\lambda \frac{t}{\varepsilon}\right) z_k(0) - ik^3 \int_0^t e^{(s-t)\lambda/\varepsilon} \rho_k(s) ds. \quad (4.2)$$

where  $\lambda = 1 - \varepsilon k^2$ . Because  $\varepsilon \in (0, 1]$  and  $k \in \mathbb{Z}$  should not be correlated,  $\lambda$  can take any value in  $(-\infty, 1)$ . For  $\lambda$  negative, this equation is unstable and cannot be solved numerically using standard tools. To overcome this, we consider the stabilized change of variable instead

$$z_k = j_k + \frac{ik}{1 + \alpha \varepsilon k^2} \rho_k$$

where  $\alpha$  is a positive constant which we shall calibrate as the study progresses. This is the same change of variable as before up to  $\mathcal{O}(\varepsilon)$ , but  $ik \rho_k$  was regularized with an elliptic operator to help with high frequencies. The problem to solve becomes

$$\begin{cases} \partial_t \rho_k = -\frac{k^2}{1 + \alpha \varepsilon k^2} \rho_k - ik z_k, \\ \partial_t z_k = -\frac{1}{\varepsilon} z_k + \frac{k^2}{1 + \alpha \varepsilon k^2} z_k - \frac{ik^3}{1 + \alpha \varepsilon k^2} \left( \alpha + \frac{1}{1 + \alpha \varepsilon k^2} \right) \rho_k. \end{cases} \quad (4.3)$$



As in (4.2), the growth of  $z_k$  is given by  $e^{-\lambda t/\varepsilon}$  if  $\lambda$  is defined by

$$\lambda = 1 - \frac{\varepsilon k^2}{1 + \alpha \varepsilon k^2} \in \left(1 - \frac{1}{\alpha}, 1\right].$$

For stability reasons  $\lambda$  must be positive, therefore we shall choose  $\alpha \geq 1$ .

Let us set  $u_k = (\rho_k, z_k)^T$  and  $\Lambda = \text{Diag}(0, 1)$  such that  $\partial_t u_k = -\frac{1}{\varepsilon} \Lambda u_k + f(u_k)$  with

$$f(u) = \begin{pmatrix} -\frac{k^2}{1 + \alpha \varepsilon k^2} u_1 - i k u_2 \\ \frac{k^2}{1 + \alpha \varepsilon k^2} u_2 - \frac{i k^3}{1 + \alpha \varepsilon k^2} \left( \alpha + \frac{1}{1 + \alpha \varepsilon k^2} \right) u_1 \end{pmatrix}. \quad (4.4)$$

In the upcoming study, we usually prefer the notation  $f(\rho, z)$  rather than  $f(u)$  so as to keep the distinction between both coordinates clear. Assuming  $|k| \leq k_{\max}$ , it is possible to bound  $f(\rho_k, z_k)$  independently of  $k$  and of  $\varepsilon$ , allowing us to apply the method developed in this paper in order to approximate every  $\rho_k$  and  $j_k$ , and eventually  $\rho(x, t)$  and  $j(x, t)$ . Note that no rigorous aspects of convergence in functional spaces are considered here – this will be treated in a forthcoming work. We omit the index  $k$  going forward for the sake of clarity.

The micro-macro method is initialized by setting the change of variable  $\Omega_\tau^{[0]}(\rho, z) = (\rho, e^{-\tau} z)^T$ . The vector field followed by the macro part  $v^{[0]}$  is  $F^{[0]}$  given by

$$F^{[0]}(\rho, z) = \widehat{k}^2 \begin{pmatrix} -\rho \\ z \end{pmatrix} \quad \text{with} \quad \widehat{k} = \frac{k}{\sqrt{1 + \alpha \varepsilon k^2}}. \quad (4.5)$$

This means that the macro variable  $v^{[0]}(t)$  is given by

$$v^{[0]}(t) = \begin{pmatrix} e^{-\widehat{k}^2 t} & 0 \\ 0 & e^{\widehat{k}^2 t} \end{pmatrix} v^{[0]}(0).$$

Notice that the growth of  $v_2^{[0]}(t)$  is in  $e^{\widehat{k}^2 t}$ , which is akin to the heat equation in reverse time. This is problematic, as it is possible for  $\widehat{k}$  to be quite big. For example with  $k = 10, \alpha = 2$  and  $\varepsilon = 10^{-2}$ , one gets  $e^{\widehat{k}^2} \approx 3 \cdot 10^{14}$ . However in order to obtain the solution of (4.1),  $u_k(t) = \Omega_{t/\varepsilon}^{[0]}(v^{[0]}(t)) + w^{[0]}(t)$ , we are only interested in  $\Omega_{t/\varepsilon}^{[0]}(v^{[0]}(t))$  for the macro part, and  $\eta_{t/\varepsilon}^{[0]}(v^{[0]}(t))$  for the micro part, which only depend on  $e^{-\frac{t}{\varepsilon} \Lambda} v^{[0]}(t)$  as can be seen in the upcoming expression of  $\eta^{[0]}$  and using  $\Omega_\tau^{[0]}(u) = e^{-\tau \Lambda} u$ . This means that the interesting quantity is

$$e^{-\frac{t}{\varepsilon} \Lambda} v^{[0]}(t) = \begin{pmatrix} e^{-\widehat{k}^2 t} & 0 \\ 0 & e^{-(1 - \varepsilon \widehat{k}^2) \frac{t}{\varepsilon}} \end{pmatrix} v^{[0]}(0). \quad (4.6)$$

Recognizing  $1 - \varepsilon \widehat{k}^2 = \lambda > 0$  in this expression, it follows that  $v_2^{[0]}$  is a decreasing function of time, therefore it is bounded uniformly for all  $t, k$  and  $\varepsilon$ . Because the exact computation of  $e^{-\frac{t}{\varepsilon} \Lambda} v^{[0]}(t)$  is readily available, it is used during implementation, leaving only  $w^{[0]}$  to be computed numerically using ERK schemes. Should the reader wish to conduct their own implementation, they should use the defect

$$\eta_\tau^{[0]}(\rho, z) = \begin{pmatrix} i k e^{-\tau} z \\ \widehat{k}^2 \left( \alpha + \frac{1}{1 + \alpha \varepsilon k^2} \right) i k \rho \end{pmatrix} = \eta_0^{[0]}(\rho, e^{-\tau} z).$$

By linearity of  $f$ , the micro variable  $w^{[0]}$  follows the differential equation

$$\partial_t w^{[0]} = -\frac{1}{\varepsilon} \Lambda w^{[0]} + f(w^{[0]}) - \eta_0^{[0]} \left( e^{-\frac{t}{\varepsilon} \Lambda} v^{[0]}(t) \right), \quad w^{[0]}(0) = 0.$$

The rescaled macro variable  $e^{-\frac{t}{\varepsilon} \Lambda} v^{[0]}(t)$  is given by relation (4.6) with initial condition  $v^{[0]}(0) = u(0) = (\rho_k(0), z_k(0))^T$ .

Extending our expansion to order 1 is not trivial either. Direct application of iterations (2.15) yields

$$\Omega_\tau^{[1]}(\rho, z) = \begin{pmatrix} \rho + \varepsilon i k e^{-\tau} z \\ e^{-\tau} z - \varepsilon \hat{k}^2 \left( \alpha + \frac{1}{1 + \alpha \varepsilon k^2} \right) i k \rho \end{pmatrix}$$

from which the vector field for the macro part is

$$F^{[1]}(\rho, z) = \hat{k}^2 \left( 1 + \varepsilon k^2 \left( \alpha + \frac{1}{1 + \alpha \varepsilon k^2} \right) \right) \begin{pmatrix} -\rho \\ z \end{pmatrix}.$$

Following the same reasoning as before, one should study the evolution of the  $z$ -component of the rescaled macro variable  $e^{-\frac{t}{\varepsilon} \Lambda} v^{[1]}(t)$ . This evolution is in  $e^{-\tilde{\lambda} t / \varepsilon}$  where  $\tilde{\lambda} = 1 - \varepsilon \hat{k}^2 \left( 1 + \varepsilon k^2 \left( \alpha + \frac{1}{1 + \alpha \varepsilon k^2} \right) \right)$ . Studying  $\tilde{\lambda}$  as a function of  $\varepsilon k^2$  in  $\mathbb{R}_+$  shows that it is negative for  $\varepsilon k^2 > 1$ , whatever the value of  $\alpha \geq 1$ .

To circumvent this, we replace  $\varepsilon$  by  $\frac{\varepsilon}{1 + \alpha \varepsilon k^2}$  in iterations (2.15). This adds terms of order  $\varepsilon^2$  in the definition of  $\Omega^{[1]}$  that do not modify any properties of the micro-macro decomposition but it regularises the problem. Specifically, we define

$$\Omega_0^{[1]}(\rho, z) = \begin{pmatrix} \rho + \frac{\varepsilon}{1 + \alpha \varepsilon k^2} i k z \\ z - \frac{\varepsilon}{1 + \alpha \varepsilon k^2} \hat{k}^2 \left( \alpha + \frac{1}{1 + \alpha \varepsilon k^2} \right) i k \rho \end{pmatrix}, \quad (4.7)$$

from which we get the vector field

$$F^{[1]}(\rho, z) = \hat{k}^2 \left( 1 + \varepsilon \hat{k}^2 \left( \alpha + \frac{1}{1 + \alpha \varepsilon k^2} \right) \right) \begin{pmatrix} -\rho \\ z \end{pmatrix}.$$

This time also, the identities  $\Omega_\tau^{[1]}(u) = \Omega_0^{[1]}(e^{-\tau \Lambda} u)$  and  $\eta_\tau^{[1]}(u) = \eta_0^{[1]}(e^{-\tau \Lambda} u)$  are satisfied, therefore the interesting variable is  $e^{-\frac{t}{\varepsilon} \Lambda} v^{[1]}(t)$ . The quantity dictating its growth is

$$\tilde{\lambda} = 1 - \varepsilon \hat{k}^2 \left( 1 + \varepsilon \hat{k}^2 \left( \alpha + \frac{1}{1 + \alpha \varepsilon k^2} \right) \right)$$

which is positive for all  $\varepsilon k^2 \in \mathbb{R}_+$  if and only if  $\alpha \geq 2$ . As with the expansion of order 0, the macro variable should be rescaled and computed exactly. The micro part  $w^{[1]}$  is given by the differential equation

$$\partial_t w^{[1]} = -\frac{1}{\varepsilon} \Lambda w^{[1]} + f(w^{[1]}) - \eta_0^{[1]} \left( e^{-\frac{t}{\varepsilon} \Lambda} v^{[1]}(t) \right), \quad w^{[1]}(0) = u_k(0) - \Omega_0^{[1]} \left( v^{[1]}(0) \right) \quad (4.8)$$

where, writing  $\hat{I} = (1 + \alpha \varepsilon k^2)^{-1}$ ,

$$\eta_\tau^{[1]}(\rho, z) = i k \cdot \varepsilon \hat{k}^2 \left( \alpha + \hat{I} \left( 2 + \varepsilon \hat{k}^2 (\alpha + \hat{I}) \right) \right) \begin{pmatrix} e^{-\tau} z \\ \hat{k}^2 (\alpha + \hat{I}) \rho \end{pmatrix} \quad (4.9)$$

$$\text{and } v^{[1]}(0) = \begin{pmatrix} \rho_k(0) - \varepsilon \widehat{I} i k z_k(0) \\ z_k(0) + \varepsilon \widehat{k}^2 (\alpha + \widehat{I}) i k \rho_k(0) \end{pmatrix}. \quad (4.10)$$

We approached the initial condition using Remark 2.7, but an exact computation of the exact initial condition  $(\Omega_0^{[1]})^{-1}(u_0)$  is possible, as the map  $u \mapsto \Omega_0^{[1]}(u)$  is linear.

**Proposition 4.1.** *Given a maximum frequency  $k_{\max} > 0$  and a scalar  $\alpha \geq 2$ , and assuming  $|k| \leq k_{\max}$ , the solution  $u_k$  of problem (4.3) can be decomposed into*

$$u_k(t) = \Omega_0^{[1]} \left( e^{-\frac{t}{\varepsilon} \Lambda} v^{[1]}(t) \right) + w^{[1]}(t)$$

where  $\Omega_0^{[1]}$  is given by (4.7) and  $w^{[1]}(t) = \mathcal{O}(\varepsilon^2)$ . The macro component  $v^{[1]}$  is given by

$$e^{-\frac{t}{\varepsilon} \Lambda} v^{[1]}(t) = \begin{pmatrix} e^{-K^{[1]} t} & 0 \\ 0 & e^{-(1-\varepsilon K^{[1]}) \frac{t}{\varepsilon}} \end{pmatrix} v^{[1]}(0)$$

with  $K^{[1]} = \widehat{k}^2 \left( 1 + \varepsilon \widehat{k}^2 \left( \alpha + \frac{1}{1+\alpha \varepsilon k^2} \right) \right)$ ,  $\widehat{k} = \frac{k}{\sqrt{1+\alpha \varepsilon k^2}}$  and  $v^{[1]}(0)$  is either  $(\Omega_0^{[k]})^{-1}(u_k(0))$  or its approximation (4.10). The micro component  $w^{[1]}$  is the solution to

$$\partial_t w^{[1]} = -\frac{1}{\varepsilon} \Lambda w^{[1]} + f(w^{[1]}) - \eta_0^{[1]} \left( e^{-\frac{t}{\varepsilon} \Lambda} v^{[1]}(t) \right), \quad w^{[1]}(0) = u_k(0) - \Omega_0^{[k]}(v^{[1]}(0))$$

with  $f$  and  $\eta_0^{[1]}$  given respectively by (4.4) and (4.9). With these definitions,  $w^{[1]}$  can be computed with a uniform error of order 2, therefore  $u_k$  can be computed with a uniform error of order 2.

The reader may notice that only a finite number of modes is considered. This is required so that there exists a bound uniform w.r.t.  $k$  and  $\varepsilon$  on the micro part of the problem (4.8) in order to apply our method. This is amenable to a CFL condition, i.e. some stiffness still exists due to the nature of the problem, but this stiffness is independent of  $\varepsilon$ . This is what we mean by uniform accuracy.

## 4.2 Relaxed conservation law

Our second test case is a hyperbolic problem for  $(t, x) \in [0, T] \times \mathbb{R}/2\pi\mathbb{Z}$ ,

$$\begin{cases} \partial_t u + \partial_x \tilde{u} = 0, \\ \partial_t \tilde{u} + \partial_x u = \frac{1}{\varepsilon} (g(u) - \tilde{u}), \end{cases} \quad (4.11)$$

with smooth initial conditions  $u(0, x)$  and  $\tilde{u}(0, x)$ . This is a stiffly relaxed conservation law, as presented in [JX95].

**Remark 4.2.** *Note that the assumption that  $\Lambda$  has integer coefficients is restrictive in this case. One may want to consider the equation on the second coordinate to be*

$$\partial_t \tilde{u} + \partial_x u = \frac{\sigma(x)}{\varepsilon} (b(x)u - \tilde{u})$$

as is done in [HS21], however this is not possible with our method.

In order to proceed, we require the following condition to be met:

$$|g'(u)| < 1 \quad (4.12)$$

This is a known stability condition when deriving asymptotic expansions for this kind of problem.

We start by discretising this system in space with  $N > 0$  points. Going forward,  $(x_j)_{j \in \mathbb{Z}/N\mathbb{Z}}$  denotes a fixed uniform discretisation of  $\mathbb{R}/2\pi\mathbb{Z}$ , of mesh size  $\Delta x := 2\pi/N$ . We define the vectors  $U = (u_j)_j, \tilde{U} = (\tilde{u}_j)_j$  and, given a vector  $V = (v_j)_j$  of size  $N$ ,  $g(V) = (g(v_j))_j$ . For simplicity,  $u_j(t)$  is the approximation of  $u(t, x_j)$ , and the same goes for  $\tilde{u}$ . We denote  $D$  the matrix of centered finite differences and  $L$  the standard discrete Laplace operator, which is to say

$$DV = \left( \frac{1}{2\Delta x} (v_{j+1} - v_{j-1}) \right)_j \quad \text{and} \quad LV = \left( \frac{1}{\Delta x^2} (v_{j+1} - 2v_j + v_{j-1}) \right)_j$$

Using an upwind scheme after diagonalising problem (4.11) yields

$$\begin{cases} \partial_t U + D\tilde{U} - \frac{\Delta x}{2} LU = 0, \\ \partial_t \tilde{U} + DU - \frac{\Delta x}{2} L\tilde{U} = \frac{1}{\varepsilon} (g(U) - \tilde{U}). \end{cases} \quad (4.13)$$

Setting  $U_1 = U$  and  $U_2 := \tilde{U} - g(U_1)$ , and neglecting the terms involving  $L$  for clarity, this problem becomes

$$\begin{cases} \partial_t U_1 = -D(U_2 + g(U_1)), \\ \partial_t U_2 = -\frac{1}{\varepsilon} U_2 + g'(U_1) D U_2 - T(U_1) \end{cases} \quad (4.14)$$

where we defined  $T(U_1) := D U_1 - g'(U_1) D g(U_1)$ . From this, our method can be applied, but precautions must be taken in order to avoid having to solve the heat equation in backwards time. Therefore we set

$$\Omega_\tau^{[1]}(U_1, U_2) = \begin{pmatrix} U_1 + \varepsilon(1 - 2\varepsilon D^2)^{-1} D U_2 \\ e^{-\tau} U_2 - \varepsilon T(U_1) \end{pmatrix}.$$

Similarly to the manipulations for the telegraph equation, we multiplied  $\varepsilon$  by  $(I_N - 2\varepsilon D^2)^{-1}$ , but this time only for the first component. Writing  $\tilde{D} = (I_N - 2\varepsilon D^2)^{-1} D$ , the associated vector field is

$$F^{[1]}(U_1, U_2) = \begin{pmatrix} -Dg(U_1) + \varepsilon D T(U_1) \\ g'(U_1) D U_2 - \varepsilon T'(U_1) \tilde{D} U_2 - \varepsilon^2 g''(U_1) (T(U_1), \tilde{D} U_2) \end{pmatrix}.$$

it is possible to obtain  $\Omega^{[0]}$  and  $F^{[0]}$  by neglecting the terms of order  $\varepsilon$  and above in the expressions above.

**Remark 4.3.** Remember that for the telegraph equation, the macro variable  $v^{[1]}(t)$  needed to be rescaled by  $e^{-t\Lambda/\varepsilon}$ . This is not the case here: In the limit  $\Delta x \rightarrow 0$ , the macro variable  $v^{[1]} = (\bar{u}_1, \bar{u}_2)^T$  is given by

$$\begin{cases} \partial_t \bar{u}_1 = -\partial_x [g(\bar{u}_1) - \varepsilon(1 - g'(\bar{u}_1)^2)\partial_x \bar{u}_1], \\ \partial_t \bar{u}_2 = g'(\bar{u}_1)\partial_x \bar{u}_2 - (1 - g'(\bar{u}_1)^2) \cdot (1 - 2\varepsilon\partial_x^2)^{-1}\varepsilon\partial_x^2 \bar{u}_2 + \varepsilon\phi^\varepsilon(\bar{u}_1, \tilde{D}\bar{u}_2) \end{cases}$$

with  $\tilde{D} = (1 - 2\varepsilon\partial_x^2)^{-1}\partial_x$  and  $\phi^\varepsilon(u_1, u_2) = g''(u_1)(2g'(u_1) - \varepsilon(1 - g'(u_1)^2)\partial_x u_1)u_2$ . The operator  $(1 - 2\varepsilon\partial_x^2)^{-1}\varepsilon\partial_x^2$  is bounded, therefore  $\bar{u}_2$  is well-defined. The equation on  $\bar{u}_1$  is a well-known result. If  $\varepsilon$  was also relaxed in the  $U_2$ -component of  $\Omega^{[1]}$ , there might be no need for condition (4.12) but the result would be different.

Because  $D^2$  is sparse, it is not too costly to compute  $(I_N - \varepsilon D^2)^{-1}$ , however the conditioning may depend on the ratio between  $\varepsilon$  and  $\Delta x$ . Indeed, studying the eigenvalues of  $D$  reveals that the eigenvalues  $(\mu_k)_{k \in \mathbb{Z}/N\mathbb{Z}}$  of  $I_N - \varepsilon D^2$  are

$$\mu_k = 1 + \frac{\varepsilon}{\Delta x^2} \sin^2 \left( 2\pi \frac{k}{N} \right) \quad (4.15)$$

meaning that for  $N$  big, the conditioning is approximately  $1 + \varepsilon/\Delta x^2$ . Therefore, for  $\varepsilon$  big and  $\Delta x$  small, this inversion can become very costly, even though the cost remains bounded independently of  $\varepsilon$ .

Obtaining the defects of order 0 and 1 from these expressions presents no difficulty. For  $\eta^{[1]}$ , we separate here the  $U_1$ -component and the  $U_2$ -component for clarity.

$$\eta_\tau^{[0]}(U_1, U_2) = \begin{pmatrix} e^{-\tau} D U_2 \\ T(U_1) \end{pmatrix},$$

$$\begin{aligned} \eta_0^{[1]}(U_1, U_2)_{U_1} &= D(g(U_1 + \varepsilon \tilde{D}W) - g(U_1)) + (D - \tilde{D})U_2 \\ &\quad + \varepsilon \tilde{D} \left( g'(U_1)DW - \varepsilon T'(U_1)\tilde{D}W - \varepsilon^2 g''(U_1)(T(U_1), \tilde{D}W) \right), \end{aligned} \quad (4.16a)$$

$$\begin{aligned} \eta_0^{[1]}(U_1, U_2)_{U_2} &= -(g'(U_1 + \varepsilon \tilde{D}U_2) - g'(U_1))DU_2 \\ &\quad + T(U_1 + \varepsilon \tilde{D}U_2) - T(U_1) - \varepsilon T'(U_1)\tilde{D}U_2 \\ &\quad + \varepsilon g'(U_1 + \varepsilon \tilde{D}U_2)DT(U_1) - \varepsilon^2 g''(U_1)(\tilde{D}U_2, T(U_1)) \\ &\quad + \varepsilon T'(U_1)(Dg(U_1) - \varepsilon T(U_1)). \end{aligned} \quad (4.16b)$$

The values of  $\eta_\tau^{[1]}(U_1, U_2)$  can be recovered using the identity

$$\eta_\tau^{[1]}(U_1, U_2) = \eta_0^{[1]}(U_1, e^{-\tau}U_2).$$

## 5 Numerical simulations

In this section we shall demonstrate our results by confirming the theoretical convergence rates of exponential Runge-Kutta (ERK) schemes from [HO05]. We also use these schemes on the original problem (1.1), thereby exhibiting the problem of order reduction.

In Subsection 5.1 we study a toy model with some non-linearity that can be found in [CCS16], for which we compute the micro-macro expansion up to order 2. In Subsection 5.2, we showcase the results of uniform convergence for the partial differential equations of Section 4. For these, the exact solution shall not take into account the error in space, i.e. it will be the solution to the discretized problem. Finally in Subsection 5.3, we discuss

### 5.1 Oscillating toy problem

We first study an "oscillating" problem presented in [CCS16] which demonstrates a possible use of the method when studying non-linear problems:

$$\begin{cases} \dot{x} = (1-z) \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix} x \\ \dot{z} = -\frac{1}{\varepsilon}z + x_1^2 x_2^2 \end{cases} \quad (5.1)$$

with initial conditions  $x_0 = (0.1, 0.7)^T$  and  $z_0 = 0.05$ , and final time  $T = 1$ . This is of the form  $\partial_t u = -\frac{1}{\varepsilon}\Lambda u + f(u)$  when setting

$$u = \begin{pmatrix} x \\ z \end{pmatrix}, \quad \Lambda = \text{Diag}(0, 0, 1) \quad \text{and} \quad f(u) = \begin{pmatrix} -(1-u_3)u_2 \\ (1-u_3)u_1 \\ (u_1 u_2)^2 \end{pmatrix}.$$

The macro part of our micro-macro decomposition is built by solving iterations on the homological equation

$$(\partial_\tau + \Lambda)\Omega_\tau^{[n+1]} = \varepsilon (f \circ \Omega_\tau^{[n]} - \partial_u \Omega_\tau^{[n]} F^{[n]}) \quad (5.2)$$

where  $F^{[n]} = \langle f \circ \Omega^{[n]} \rangle$  with  $\langle \cdot \rangle$  the projector on the  $e^{-\tau\Lambda}$ -component parallel to the other components of the exponential series. We choose the initial condition  $\Omega_\tau^{[0]} = e^{-\tau\Lambda}$  and closure condition  $\langle \Omega^{[0]} \rangle = e^{-\tau\Lambda}$ . The first iteration yields

$$\Omega_\tau^{[1]}(x, z) = \begin{pmatrix} x_1 - \varepsilon e^{-\tau} x_2 z \\ x_2 + \varepsilon e^{-\tau} x_1 z \\ e^{-\tau} z + \varepsilon (x_1 x_2)^2 \end{pmatrix} \quad \text{and} \quad F^{[1]}(x, z) = \begin{pmatrix} -(1 - \varepsilon (x_1 x_2)^2) x_2 \\ (1 - \varepsilon (x_1 x_2)^2) x_1 \\ 2\varepsilon x_1 x_2 z (x_1^2 - x_2^2) \end{pmatrix}.$$

In order to compute the second order decomposition, one must compute the difference  $T^{[1]} = f \circ \Omega^{[1]} - \partial_u \Omega^{[1]} F^{[1]}$ , which is also used to compute the defect  $\delta^{[1]} = \frac{1}{\varepsilon}(\partial_\tau + \Lambda)\Omega^{[1]} - T^{[1]}$ . From a direct calculation this writes,

$$T_\tau^{[1]}(x, z) = \begin{pmatrix} e^{-\tau} z (x_2 + \varepsilon e^{-\tau} x_1 z + 2\varepsilon^2 x_1 x_2^2 (x_1^2 - x_2^2)) \\ -e^{-\tau} z (x_1 - \varepsilon e^{-\tau} x_2 z - 2\varepsilon^2 u_1^2 u_2 (x_1^2 - x_2^2)) \\ Z_0 + \varepsilon Z_1 + \varepsilon^2 Z_2 \end{pmatrix}$$

where for clarity we defined

$$Z_0 = (x_1^2 + \varepsilon^2 e^{-2\tau} (x_2 z)^2) (x_2 + \varepsilon^2 e^{-2\tau} (x_1 z)^2),$$

$$Z_1 = -2x_1 x_2 (x_1^2 - x_2^2) (1 - \varepsilon (x_1 x_2)^2 + \varepsilon e^{-3\tau} z^3) \quad \text{and} \quad Z_2 = -e^{-2\tau} (2u_1 u_2 u_3)^2.$$

To compute the expansion of order 2, we truncate terms of order  $\varepsilon^2$  and above in  $T^{[1]}$  (which will not impact results of uniform accuracy) and solve (5.2). This yields<sup>4</sup>

$$\Omega_\tau^{[2]}(x, z) = \begin{pmatrix} x_1 - \varepsilon e^{-\tau} x_2 z - \frac{1}{2} \varepsilon^2 e^{-2\tau} z^2 x_1 \\ x_2 + \varepsilon e^{-\tau} x_1 z - \frac{1}{2} \varepsilon^2 e^{-2\tau} z^2 x_2 \\ z + \varepsilon (x_1 x_2)^2 - 2\varepsilon^2 x_1 x_2 (x_1^2 - x_2^2) \end{pmatrix},$$

$$F^{[2]}(x, z) = \begin{pmatrix} x_2(-1 + \varepsilon (x_1 x_2)^2 - 2\varepsilon^2 x_1 x_2 (x_1^2 - x_2^2)) \\ x_1(1 - \varepsilon (x_1 x_2)^2 + 2\varepsilon^2 x_1 x_2 (x_1^2 - x_2^2)) \\ 2\varepsilon z x_1 x_2 (x_1^2 - x_2^2) \end{pmatrix}.$$

The defect  $\eta^{[2]}$  is obtained using relation (2.17) or by computing  $\delta^{[2]}$  and identifying the Fourier coefficients.

**Remark 5.1.** *It is possible to find an approximation of the center manifold  $x \mapsto \varepsilon h^\varepsilon(x)$  by taking the limit  $\tau \rightarrow \infty$  of the  $z$ -component of  $\Omega^{[k]}$ . For example here*

$$\varepsilon h^\varepsilon(x) = \varepsilon (x_1 x_2)^2 - 2\varepsilon^2 x_1 x_2 (x_1^2 - x_2^2) + \mathcal{O}(\varepsilon^3).$$

*This coincides with the results in [CCS16].*

We remind the reader that the problem that is solved at times  $(t_i)_{0 \leq i \leq N}$  is

$$\begin{cases} \partial_t v^{[k]}(t) = F^{[k]}(v^{[k]}), \\ \partial_t w^{[k]}(t) = -\frac{1}{\varepsilon} \Lambda w^{[k]} + f(\Omega_{t/\varepsilon}^{[k]}(v^{[k]}) + w^{[k]}) - f(\Omega_{t/\varepsilon}^{[k]}(v^{[k]})) - \eta_{t/\varepsilon}^{[k]}(v^{[k]}), \end{cases}$$

with  $k = 1, 2$ . This yields vectors  $(v_i) \approx (v^{[k]}(t_i))$  and  $(w_i) \approx (w^{[k]}(t_i))$ , from which an approximation  $u_i \approx u^\varepsilon(t_i)$  is then obtained by setting  $u_i = \Omega_{t_i/\varepsilon}^{[k]}(v_i) + w_i$ . Initial conditions  $v^{[k]}(0)$  and  $w^{[k]}(0)$  are computed using Remark 2.7.

The difference  $f(\Omega_{t/\varepsilon}^{[2]}(v^{[2]}) + w^{[2]}) - f(\Omega_{t/\varepsilon}^{[2]}(v^{[2]}))$  is computed using

$$f(x + \tilde{x}, z + \tilde{z}) - f(x, z) = \begin{pmatrix} -(1-z)\tilde{x}_2 + (x_2 + \tilde{x}_2)\tilde{z} \\ (1-z)\tilde{x}_1 - (x_1 + \tilde{x}_1)\tilde{z} \\ (x_1 x_2 + (x_1 + \tilde{x}_1)(x_2 + \tilde{x}_2))(x_1 \tilde{x}_2 + \tilde{x}_1 x_2 + \tilde{x}_1 \tilde{x}_2) \end{pmatrix}$$

in order to avoid rounding errors due to the size difference between  $u$  and  $\tilde{u}$ .

Figure 1 showcases the phenomenon of order reduction when solving the original problem (5.1): Despite using a scheme of order 2, the error depends of  $\varepsilon$  in such a way that there exists no constant  $C$  such that the error is bounded by  $C\Delta t^2$  for all  $\varepsilon$ . However there exists  $C$  such that the error is bounded by  $C\Delta t$ . In that case, we cannot say that the error is of *uniform* order 2, as this would require the error to be independent of  $\varepsilon$ . However, this is the case when solving the micro-macro problem, as can be seen on the right-hand side of Figure 1 for a decomposition of order 2. Furthermore, the theoretical orders of convergence from Theorem 2.8 are confirmed. Indeed, using a scheme of order 2 (resp. 3) on the micro-macro problem of order 1 (resp. 2) generates a uniform error of the expected order of convergence, with no order reduction.

<sup>4</sup>It has been pointed out to the authors that the same result is obtained using nonlinear coordinate transforms described in [Rob14]. Some normal form methods compiled in [Mur06] also yield this result.

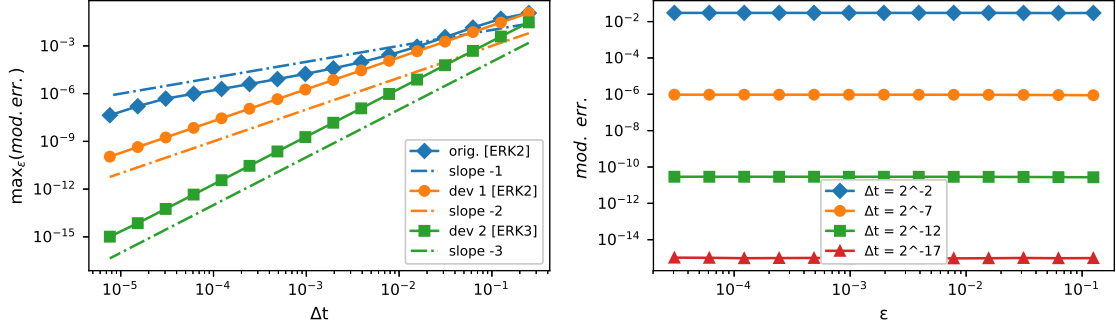


Figure 1: *Oscillating case: On the left, maximum error on  $\varepsilon$  (for  $\varepsilon = 2^{-k}$  with  $k$  spanning  $\{3, \dots, 15\}$ ) as a function of  $\Delta t$  when using exponential RK schemes (abbr. ERK) of different orders. On the right, the error as a function of  $\varepsilon$  when solving the micro-macro problem of order 2 using ERK3.*

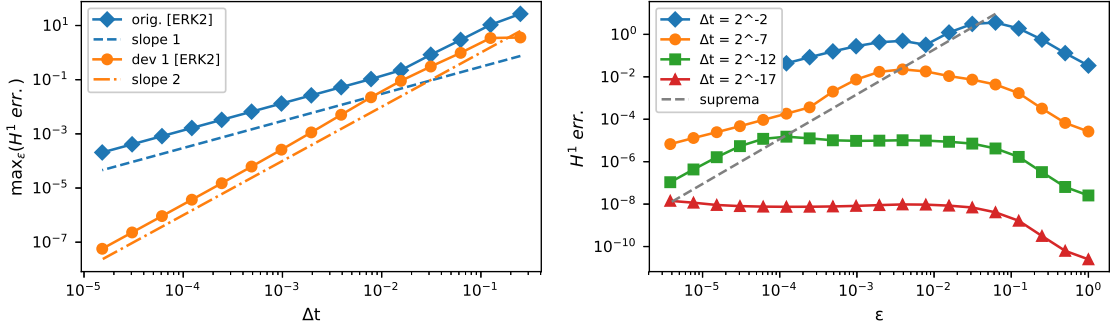


Figure 2: *Telegraph equation: Absolute  $H^1$  error on the solution of (4.1) computed by an ERK3 scheme. Supremum on  $\varepsilon$  as a function of  $\Delta t$  (left) and evolution of this error as a function of  $\varepsilon$  for the 1st-order decomposition (right).*

## 5.2 Discretized hyperbolic partial differential equations

### The telegraph equation

Using a spectral decomposition, we solve the problem, for  $(t, x) \in [0, T] \times \mathbb{R}/2\pi\mathbb{Z}$ ,

$$\begin{cases} \partial_t \rho + \partial_x j = 0, \\ \partial_t j + \frac{1}{\varepsilon} \partial_x \rho = -\frac{1}{\varepsilon} j, \end{cases}$$

by setting  $z = j + (1 - \alpha\varepsilon\Delta)^{-1} \partial_x z$ , yielding problem (4.3). The micro-macro decomposition of order 1 is summarized in Property 4.1, and its construction is detailed in Subsection 4.1. Implementations are conducted using  $\alpha = 2$ , space frequencies are bounded by  $k_{\max} := 12$ , and initial data is  $\rho(0, x) = e^{\cos(x)}$ ,  $j(0, x) = \frac{1}{2} \cos^3(x)$ .

Results can be seen in Figure 2 when using a scheme of order 2. When solving the original problem, the uniform order degenerates from 2 to 1. When considering the micro-macro problem, the order of convergence is not reduced and stays of order 2. Although it varies with  $\varepsilon$  when considering a fixed  $\Delta t$ , when considering the supremum on  $\varepsilon$ , there is no order reduction. The dashed slope on the right plot interpolates the position of the supremum of the error for each fixed  $\Delta t$ . While the error seems to improve for  $\varepsilon \ll \Delta t$ , this does not cause any order reduction. This is stronger than the property of preservation of asymptotes (which ERK schemes have, see [DP11]), since AP schemes only need to be well-defined in



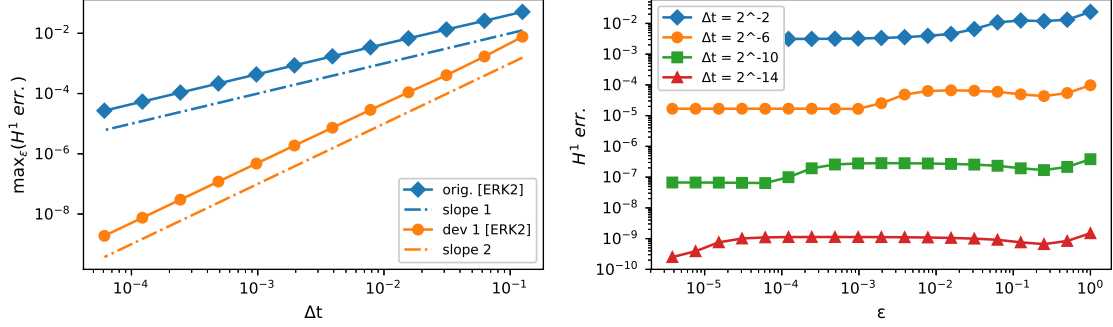


Figure 3: *Relaxed Burgers-type problem: Maximum modified  $H^1$  error (for  $\varepsilon$  spanning 1 to  $2^{-18}$  using an ERK3 scheme as a function of  $\Delta t$  (left), and  $H^1$  error as a function of  $\varepsilon$  for the micro-macro problem of order 1 (right).*

the limit  $\varepsilon \rightarrow 0$ . For them, this supremum does not need to be bounded. It appears that the relationship between the error bound and the stiffness of the linear operator is rather complex when using exponential RK schemes (again, see [HO05] for details).

### Relaxed conservation law

Our second test case is a hyperbolic problem for  $(t, x) \in [0, T] \times \mathbb{R}/2\pi\mathbb{Z}$ ,

$$\begin{cases} \partial_t u + \partial_x \tilde{u} = 0, \\ \partial_t \tilde{u} + \partial_x u = \frac{1}{\varepsilon}(g(u) - \tilde{u}), \end{cases}$$

discretized with finite volumes and written in the form of (1.1) by setting  $u_1 = u$  and  $u_2 = \tilde{u} - g(u)$  the  $x^\varepsilon$ - and the  $z^\varepsilon$ -component respectively. The micro-macro expansion is computed to order 1 using the strategy detailed in Subsection 4.2.

For our tests, following [HS21], we consider  $g(u) = bu^2$  with  $b = 0.2$ . Simulations run to a final time  $T = 0.25$  and the mesh size is fixed:  $N = 16$ . Initial data is  $u(0, x) = \frac{1}{2}e^{\sin(x)}$  and  $\tilde{u}(0, x) = \cos(x)$ . The reference solution was computed up to a precision  $10^{-12}$  using an ERK2 scheme. Convergence results are presented in Figure 3, confirming theoretical results once more.

It should be said again that our approach does not study the error in space, only in time. For instance, the relationship between the error bound and the grid size is not considered. Further studies will be conducted, especially considering CFL conditions,  $L^2$  and  $H^1$  norms, and computational costs.

## 5.3 Thoughts

### Computing cost

Note that when using a given scheme, solving a single step is much more costly for the micro-macro problem than for the direct problem: Not only is the system size doubled, but the functions implicated require more computing power to obtain a single value (especially the defect, see (4.16) for instance). It is therefore plausible to think that our method is best for computing values during the transient phase, after which it is possible to solve the original problem with uniform accuracy.

The regularized derivation  $(I_N - 2\varepsilon D^2)^{-1}D$  which appears in the micro-macro problem of the relaxed hyperbolic system may be prohibitively costly to compute for some schemes

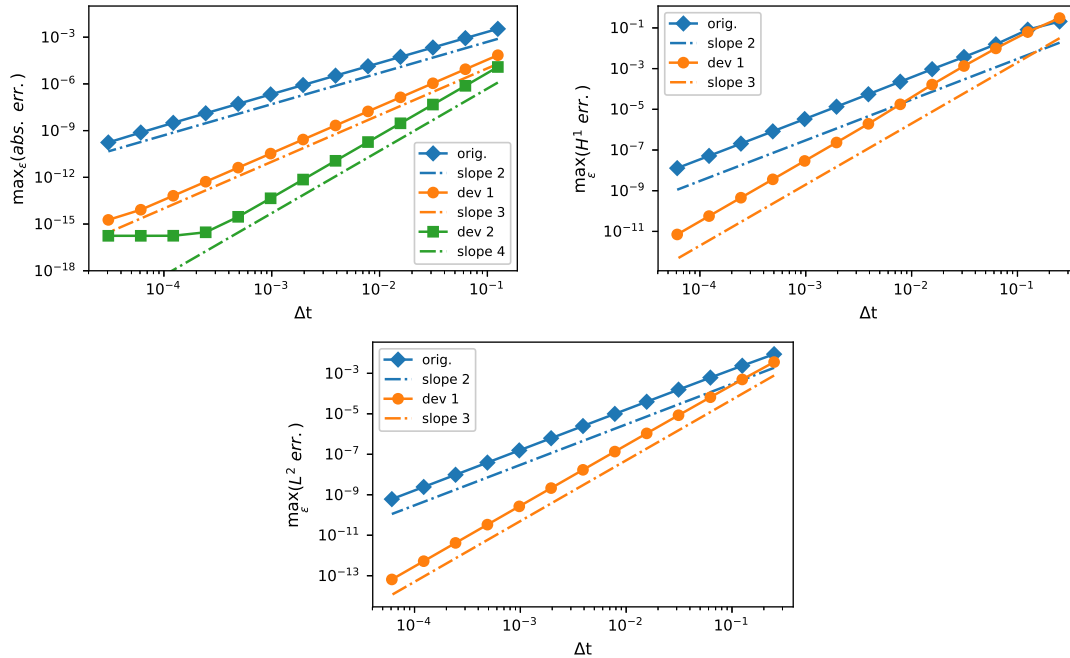


Figure 4: *In reading order, errors when solving the oscillating toy problem, the telegraph equation and the relaxed conservation law. All systems start near equilibrium and are solved with exponential Runge-Kutta schemes of the observed order of convergence.*

such as WENO, for which the derivation operator is non-linear. However we may be able to work around this, as the goal of the relaxation term is only to dampen high-frequencies, and as such inverting any discrete Laplace operator should suffice, independently of the scheme used to discretize the transport. Clearly, the subject of utilizing such regularizations for numerical purposes is complex and beyond the scope of this paper.

### Near-equilibrium convergence

If one chooses an initial condition  $z^\varepsilon(0) = 0$  in (1.1), then it is close to the center manifold up to  $\mathcal{O}(\varepsilon)$ , and Problem (1.2) can be solved with uniform accuracy of order 2 but only when considering the absolute error  $|\cdot|$ , not the modified error  $|\cdot|_\varepsilon$  from (2.22). The same behaviour is observed for the telegraph equation when setting  $j(0, x) = -\partial_x \rho(0, x)$ , meaning  $z = \mathcal{O}(\varepsilon)$ . This would theoretically mean that we need to push the micro-macro decompositions up to order 2 if we want to improve the order of convergence. However, this is not the case: uniform accuracy of order 3 is obtained from an expansion of order 1 for all test cases. This "order gain" also propagates to our micro-macro decomposition of order 2 for the oscillating toy problem. These results can be seen in Figure 4 and will be studied in future works.

### Acknowledgements

This work has been carried out within the framework of the EUROfusion Consortium and has received funding from the Euratom research and training programme 2014- 2018 and 2019-2020 under grant agreement No 633053. The views and opinions expressed herein do not necessarily reflect those of the European Commission.

## References

- [ACM99] Georgios Akrivis, Michel Crouzeix, and Charalambos Makridakis. “Implicit-explicit multistep methods for quasilinear parabolic equations”. In: *Numerische Mathematik* 82.4 (1999), pp. 521–541.
- [ADP20] Giacomo Albi, Giacomo Dimarco, and Lorenzo Pareschi. “Implicit-explicit multistep methods for hyperbolic systems with multiscale relaxation”. In: *SIAM Journal on Scientific Computing* 42.4 (2020), A2402–A2435.
- [ARW95] Uri M Ascher, Steven J Ruuth, and Brian TR Wetton. “Implicit-explicit methods for time-dependent partial differential equations”. In: *SIAM Journal on Numerical Analysis* 32.3 (1995), pp. 797–823.
- [AP96] Pierre Auger and Jean-Christophe Poggiale. “Emergence of population growth models: fast migration and slow growth”. In: *Journal of Theoretical Biology* 182.2 (1996), pp. 99–108.
- [BPR17] Sebastiano Boscarino, Lorenzo Pareschi, and Giovanni Russo. “A unified IMEX Runge–Kutta approach for hyperbolic systems with multiscale relaxation”. In: *SIAM Journal on Numerical Analysis* 55.4 (2017), pp. 2085–2109.
- [Car82] Jack Carr. *Applications of centre manifold theory*. Vol. 35. Applied Mathematical Sciences. Springer-Verlag New York, 1982.
- [CCS18] Francois Castella, Philippe Chartier, and Julie Sauzeau. “Analysis of a time-dependent problem of mixed migration and population dynamics”. In: *arXiv preprint, arXiv:1512.01880* (2018).
- [CCMM15] François Castella, Philippe Chartier, Florian Méhats, and Ander Murua. “Stroboscopic Averaging for the Nonlinear Schrödinger Equation”. In: *Foundations of Computational Mathematics* 15.2 (Apr. 2015), pp. 519–559.
- [CCS16] François Castella, Philippe Chartier, and Julie Sauzeau. “A formal series approach to the center manifold theorem”. In: *Foundations of Computational Mathematics* (2016), pp. 1–38.
- [CLMV19] Philippe Chartier, Mohammed Lemou, Florian Méhats, and Gilles Vilmart. “A New Class of Uniformly Accurate Numerical Schemes for Highly Oscillatory Evolution Equations”. In: *Foundations of Computational Mathematics* (2019).
- [CLMZ20] Philippe Chartier, Mohammed Lemou, Florian Méhats, and Xiaofei Zhao. “Derivative-free high-order uniformly accurate schemes for highly-oscillatory systems”. In: *submitted preprint* (2020).
- [DP11] Giacomo Dimarco and Lorenzo Pareschi. “Exponential Runge–Kutta methods for stiff kinetic equations”. In: *SIAM Journal on Numerical Analysis* 49.5 (2011), pp. 2057–2077.
- [GHM94] Günther Greiner, JAP Heesterbeek, and Johan AJ Metz. “A singular perturbation theorem for evolution equations and time-scale arguments for structured population models”. In: *Canadian applied mathematics quarterly* 3.4 (1994), pp. 435–459.
- [HW96] Ernst Hairer and Gerhard Wanner. *Solving ordinary differential equations II. Stiff and Differential-Algebraic Problems*. Springer Berlin Heidelberg, 1996.
- [HO04] Marlis Hochbruck and Alexander Ostermann. “Exponential Runge–Kutta methods for parabolic problems”. In: *Applied Numerical Mathematics* 53.2-4 (2004), pp. 323–339.
- [HO05] Marlis Hochbruck and Alexander Ostermann. “Explicit exponential Runge–Kutta methods for semilinear parabolic problems”. In: *SIAM Journal on Numerical Analysis* 43.3 (2005), pp. 1069–1090.

- [HS21] Jingwei Hu and Ruiwen Shu. “On the uniform accuracy of implicit-explicit backward differentiation formulas (IMEX-BDF) for stiff hyperbolic relaxation systems and kinetic equations”. In: *Mathematics of Computation* 90.328 (2021), pp. 641–670.
- [HR07] Willem Hundsdorfer and Steven J Ruuth. “IMEX extensions of linear multistep methods with general monotonicity and boundedness properties”. In: *Journal of Computational Physics* 225.2 (2007), pp. 2016–2042.
- [Jin99] Shi Jin. “Efficient asymptotic-preserving (AP) schemes for some multiscale kinetic equations”. In: *SIAM Journal on Scientific Computing* 21.2 (1999), pp. 441–454.
- [JPT98] Shi Jin, Lorenzo Pareschi, and Giuseppe Toscani. “Diffusive relaxation schemes for multiscale discrete-velocity kinetic equations”. In: *SIAM Journal on Numerical Analysis* 35.6 (1998), pp. 2405–2439.
- [JPT00] Shi Jin, Lorenzo Pareschi, and Giuseppe Toscani. “Uniformly accurate diffusive relaxation schemes for multiscale transport equations”. In: *SIAM Journal on Numerical Analysis* 38.3 (2000), pp. 913–936.
- [JX95] Shi Jin and Zhouping Xin. “The relaxation schemes for systems of conservation laws in arbitrary space dimensions”. In: *Communications on pure and applied mathematics* 48.3 (1995), pp. 235–276.
- [LM08] Mohammed Lemou and Luc Mieussens. “A new asymptotic preserving scheme based on micro-macro formulation for linear kinetic equations in the diffusion limit”. In: *SIAM Journal on Scientific Computing* 31.1 (2008), pp. 334–368.
- [MZ09] Stefano Maset and Marino Zennaro. “Unconditional stability of explicit exponential Runge-Kutta methods for semi-linear ordinary differential equations”. In: *Mathematics of computation* 78.266 (2009), pp. 957–967.
- [Mur06] James Murdock. *Normal forms and unfoldings for local dynamical systems*. Springer Science & Business Media, 2006.
- [Rob14] Anthony John Roberts. *Model emergent dynamics in complex systems*. Vol. 20. SIAM, 2014. Chap. IV.
- [Sak90] Kunimochi Sakamoto. “Invariant manifolds in singular perturbation problems for ordinary differential equations”. In: *Proceedings of the Royal Society of Edinburgh Section A: Mathematics* 116.1-2 (1990), pp. 45–78.
- [SAAP00] Eva Sánchez, Ovide Arino, Pierre Auger, and Rafael Bravo de la Parra. “A singular perturbation in an age-structured population model”. In: *SIAM Journal on Applied Mathematics* 60.2 (2000), pp. 408–436.
- [Vas63] Adelaida Borisovna Vasil’eva. “Asymptotic behaviour of solutions to certain problems involving non-linear differential equations containing a small parameter multiplying the highest derivatives”. In: *Russian Mathematical Surveys* 18.3 (1963), p. 13.